

# REAL-WORLD DATA AND CALIBRATED SIMULATION SUITE FOR OFFLINE TRAINING OF REINFORCEMENT LEARNING AGENTS TO OPTIMIZE ENERGY AND EMISSION IN BUILDINGS FOR ENVIRONMENTAL SUSTAINABILITY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Commercial office buildings contribute 17 percent of Carbon Emissions in the US, according to the US Energy Information Administration (EIA), and improving their efficiency will reduce their environmental burden and operating cost. A major contributor of energy consumption in these buildings are the Heating, Ventilation, and Air Conditioning (HVAC) devices. HVAC devices form a complex and interconnected thermodynamic system with the building and outside weather conditions, and current setpoint control policies are not fully optimized for minimizing energy use and carbon emission. Given a suitable training environment, a Reinforcement Learning (RL) agent is able to improve upon these policies, but training such a model, especially in a way that scales to thousands of buildings, presents many practical challenges. Most existing work on applying RL to this important task either makes use of proprietary data, or focuses on expensive and proprietary simulations that may not be grounded in the real world. We present the Smart Buildings Control Suite, the first open source interactive HVAC control dataset extracted from live sensor measurements of devices in real office buildings. The dataset consists of two components: six years of real-world historical data from three buildings, for offline RL, and a lightweight interactive simulator for each of these buildings, calibrated using the historical data, for online and model-based RL. For ease of use, our RL environments are all compatible with the OpenAI gym environment standard. We also demonstrate a novel method of calibrating the simulator, as well as baseline results on training an RL agent on the simulator, predicting real-world data, and training an RL agent directly from data. We believe this benchmark will accelerate progress and collaboration on building optimization and environmental sustainability research.

## 1 INTRODUCTION

Energy optimization and management in commercial buildings is a very important problem, whose importance is only growing with time. Buildings account for 37% of all US carbon emissions, with commercial buildings alone taking up a staggering 17% in 2023 (EIA). Reducing those emissions by even a small percentage can have a significant effect. In climates that are either very hot or very cold, energy consumption is much higher, and there is even more room to have a major impact. We believe this problem is one of the most important avenues for climate sustainability research, where even a small improvement over baseline policies can drastically reduce our carbon footprint.

In particular, HVAC systems account for 40-60% of energy use in buildings (Pérez-Lombard et al., 2008), and roughly 15% of the world’s total energy consumption (Asim et al., 2022). Most office buildings are equipped with advanced HVAC devices, like Variable Air Volume (VAV) devices, Hot Water Systems, Air Conditioners and Air Handlers that are configured and tuned by the engineers, manufacturers, installers, and operators to run efficiently with the device’s local control loops (McQuiston et al., 2023). However, integrating multiple HVAC devices from diverse vendors into a

054 building “system” requires technicians to program fixed operating conditions for these units, which  
055 may not be optimal for every building and every potential weather condition. Existing setpoint control  
056 policies are not optimal under all conditions, and the possibility exists that a machine learning  
057 model may be trained to continuously tune a small number of setpoints to achieve greater energy  
058 efficiency and reduced carbon emission.

059 Optimizing HVAC control has been an active research area for decades, and yet while AI has begun  
060 to revolutionize many industries, to date almost all HVAC systems remain the same as they were 30  
061 years ago: despite all the literature on the topic, there is not a single solution that has been widely  
062 adopted in the real world.

063 One of the most significant factors limiting progress is the lack of a reliable public benchmark to  
064 test solutions against. Current work generally makes use of proprietary data and expensive (often  
065 also proprietary) simulations. This limits participation to those with exclusive access, and makes  
066 most claims difficult to verify and compare. A strong public dataset would facilitate collaborations  
067 between institutions, standardize research efforts, and allow for wider participation. Historically,  
068 much of progress in AI has been driven by easily accessible public benchmarks, from the ImageNet  
069 Challenge in Vision (Russakovsky et al., 2015), to the Atari57 suite in RL (Badia et al., 2020), and  
070 the GLUE Benchmark in language (Wang et al., 2018). A similar benchmark in HVAC control may  
071 help accelerate progress and finally lead to adoption of solutions in the real world.

072 We present The Smart Buildings Control Suite, a high quality, fully accessible, building control  
073 benchmark. The benchmark consists of two components:

- 075 • Real-world historical HVAC data, collected from three buildings over a six year period.
- 076 • A highly customizable and scalable HVAC and building simulator, with configurations  
077 corresponding to each of the above buildings

078  
079 Our contributions include one of the first public real-world HVAC datasets, a highly customizable  
080 and scalable HVAC and building simulator, a rapid configuration method to customize the simulator  
081 to a particular building, a calibration method to improve this fidelity using real-world data, and an  
082 evaluation method to measure the simulator fidelity. The dataset contains information from three  
083 buildings in California, the largest of which is three stories and 118,086 ft<sup>2</sup>. Using data we obtained  
084 from each building, we calibrate our simulator, and demonstrate using our evaluation pipeline that  
085 this significantly improves its fidelity to the real building. We provide pre-calibrated simulators for  
086 all of our buildings, as well as code to both reproduce the calibration procedure, and to calibrate the  
087 simulator to new scenarios. While our suite focuses on three buildings, our simulator is easily adapt-  
088 able, allowing for the development of general purpose solutions that can be applied to any building.  
089 All the data and simulator code is open source and compatible with the OpenAI gym environment  
090 standard (Brockman et al., 2016), and data is available on the popular TensorFlow Datasets platform  
091 (TFDS) under the Creative Commons License.

092 We first give an overview of the problem and related work, and then present the structure of the  
093 data. Next we introduce the simulator, and discuss our configuration, calibration, and evaluation  
094 techniques. After that, we run through an example of the process of calibrating the simulator to  
095 real data, and finally we demonstrate success on three key benchmark tasks: training an RL agent  
096 on the calibrated simulator environment using Soft Actor Critic (Haarnoja et al., 2018), training a  
097 regression model to predict the real world dynamics, and training a Soft Actor Critic agent from the  
098 real world data via the regression model.

## 099 2 OPTIMIZING ENERGY AND EMISSION IN OFFICE BUILDINGS WITH RL

100  
101 In this section we frame energy optimization in office buildings as an RL problem. We define the  
102 state of the office building  $S_t$  at time  $t$  as a fixed length vector of measurements from sensors on  
103 the building’s devices, such as a specific VAV’s zone air temperature, gas meter’s flow rate, etc. The  
104 action on the building  $A_t$  is a fixed-length vector of device setpoints selected by the agent at time  $t$ ,  
105 such as the boiler supply water temperature setpoint, etc.

106  
107 More generally, RL is a branch of machine learning that attempts to train an agent to choose the best  
actions to maximize some long-term, cumulative reward (Sutton & Barto, 2018). The agent observes

the state  $S_t$  from the environment at time  $t$ , then chooses action  $A_t$ . The environment responds by transitioning to the next state  $S_{t+1}$  and returns a reward (or penalty) after the action,  $R_{t+1}$ . Over time, the agent will explore the action space and learn to maximize the reward over the long term for each given state. A discount factor  $\gamma$  reduces the value of future rewards amplifying the value of the near-term reward. When this cycle is repeated over multiple episodes, the agent converges on a state-action policy that maximizes the long-term reward.

This sequence is often formalized as the Markov Decision Process (MDP) (Garcia & Rachelson, 2013), described by the tuple  $(S, A, p, R)$  where the state space is continuous (e.g., temperatures, flow rates, etc.) and the action space is continuous (e.g., setpoint temperatures) and the transition probability  $p : S \times S \times A \rightarrow [0, 1]$  represents the probability density of the next state  $S_{t+1}$  from taking action  $A_t$  on the current state  $S_t$ . The reward function  $R : S \times A \rightarrow [R_{min}, R_{max}]$  emits a single scalar value at each time  $t$ . The agent is acting under a policy  $\pi_\theta(A_t|S_t)$  parameterized by  $\theta$  that represents the probability of taking action  $A_t$  from state  $S_t$ . The goal of an RL agent is to find the policy that maximizes the expected long-term cumulative, discounted reward. The set of parameters  $\theta^*$  of the optimal policy can be expressed as:

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \sum_t \gamma^t R(S_t, A_t) \right]$$

where  $\theta$  is the current policy parameter, and  $\tau$  is a trajectory of states, actions, and rewards over sequential time steps  $t$ . In order to converge to the optimal policy, the agent requires many training iterations to explore the policy space, making online training directly on the real-world building from scratch inefficient, dangerous, impracticable, and likely impossible. Therefore, it is necessary to enable offline learning, where the agent can train in an efficient sandbox environment that adequately emulates the dynamics of the building before being deployed to the real world.

**Reward Function** RL generally requires a single scalar reward signal,  $R_t(S_t, A_t)$  that indicates the quality of taking action  $A_t$  in state  $S_t$ . We thus define a custom feedback signal,  $R_{3C}$ , as a weighted sum of negative cost functions for carbon emission, energy cost, and comfort levels within the building, which we dubbed the 3C Reward. It is governed by the following equation:

$$R_{3C} = u \times C_1 + v \times C_2 + w \times C_3$$

where  $C_1$  represents normalized comfort conditions,  $C_2$  normalized energy cost and  $C_3$  normalized carbon emission. Constants  $u, v, w$  represent operator preferences, allowing them to weight the relative importance of cost, comfort and carbon consumption.  $R_{3C} = 0$  when no energy is consumed, no carbon is emitted, and all occupied zones are in setpoint bounds, and negative otherwise. For more details, and equations governing how we normalize and measure these quantities, see Appendix A.

### 3 RELATED WORKS

Considerable attention has been paid to HVAC control (Fong et al., 2006) in recent years (Kim et al., 2022), and while alternative approaches exist, such as model predictive control (Taheri et al., 2022), a growing portion of the literature has considered how RL and its various associated algorithms can be leveraged (Yu et al., 2021; Mason & Grijalva, 2019; Yu et al., 2020; Gao & Wang, 2023; Wang et al., 2023; Vázquez-Canteli & Nagy, 2019; Zhang et al., 2019b; Fang et al., 2022; Zhang et al., 2019b). As mentioned above, a central requirement in RL is the offline environment that trains the RL agent. Several methods have been proposed, largely falling under three broad categories.

**Data-driven Emulators** Some works attempt to learn a dynamics as a multivariate regression model from real-world data (Zou et al., 2020; Zhang et al., 2019a), often using recurrent neural network architecture, such as Long Short-Term Memory (LSTM) (Velswamy et al., 2017; Sendra-Arranz & Gutiérrez, 2020; Zhuang et al., 2023). The difficulty here is that data-driven models often do not generalize well to circumstances outside the training distribution, especially since they are not physics based.

**Offline RL** The second approach is to train the agent directly from the historical real-world data, without ever producing an interactive environment (Chen et al., 2020; 2023; Blad et al., 2022). While the real-world data is obviously of high accuracy and quality, this presents a major challenge, since the agent cannot take actions in the real world and interact with any form of an environment. This inability to explore severely limits its ability to improve over the baseline policy producing the real-world data (Levine et al., 2020). Furthermore, prior to our work, there are few public datasets available.

**Physics-based Simulation** HVAC system simulation has long been studied (Trčka & Hensen, 2010; Riederer, 2005; Park et al., 1985; Trčka et al., 2009; Husaunndee et al., 1997; Trcka et al., 2007; Blonsky et al., 2021). EnergyPlus (Crawley et al., 2001), a high-fidelity simulator developed by the Department of Energy, is commonly used (Wei et al., 2017; Azuatalam et al., 2020; Zhao et al., 2015; Wani et al., 2019; Basarkar, 2011), but suffers from scalability and configuration challenges.

To overcome the limitations of each of the above three methods, some work has proposed a hybrid approach (Zhao et al., 2021; Balali et al., 2023; Goldfeder & Sipple, 2023; Zhang et al., 2023; Klanatsky et al., 2023; Drgoňa et al., 2021), and indeed this is the category our work falls under. What is unique about our approach is the use of a physics based simulator that achieves an ideal balance between speed of configuration, and fidelity to the real world. Our simulator is lightweight enough to be configured to an arbitrary building in a matter of hours, and using our calibration process based on real-world data, accurate enough to train an effective control agent off-line. This allows our solution to be highly scalable, like the first two approaches, but still rooted in physics, and demonstrably calibrated, like the third approach.

Various works have also discussed how exactly to apply RL to an HVAC environment, such as what sort of agent to train. Inspired by prior effective use of Soft Actor Critic (SAC) on related problems (Kathirgamanathan et al., 2021; Coraci et al., 2021; Campos et al., 2022; Biemann et al., 2021), we chose to demo our environment using a SAC agent.

**Prior Datasets** While many building datasets exist (Ye et al., 2019), most either have a different focus (Sachs et al., 2012; Urban et al.; Kriechbaumer & Jacobsen, 2018; Granderson et al., 2023), do not contain sufficient HVAC information (Miller et al., 2020; Mathew et al., 2015; Rashid et al., 2019; Jazizadeh et al., 2018; Sartori et al., 2023), are focused on residential buildings (Murray et al., 2017; Barker et al., 2012; Meinrenken et al., 2020) or non-standard buildings (Pettit et al., 2014; Naug & Chandan), or are simulated (Field et al., 2010; Bakker et al., 2022). Even the few datasets directly relevant (Luo et al., 2022; Heer et al., 2024) are non-interactive. As far as we are aware, we present the first HVAC control benchmark that has high quality real-world data with computationally cheap simulations of the same buildings, allowing for both real-world grounding and interactive control experiments.

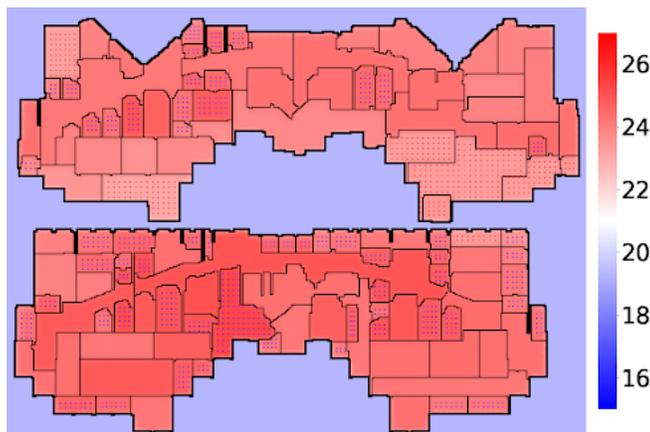


Figure 1: Example Visualization of an Environment. Blue represents colder temperatures, red warmer. Blue and red dots inside the building indicate diffusers that are dispensing cold and warm air respectively.

## 4 THE DATASET STRUCTURE

Both the real-world data and simulated data are given in the same format. Following the RL paradigm, data is provided as a series of observations, actions, and rewards. In the case of the

real-world data, this comes in the form of static historical episodes, where the actions follow the baseline policy in the building, and in the case of the simulator, as a proper interactive RL environment where actions can be taken in real time.

To make the task as realistic as possible, we formatted the data to closely resemble the real-world building API, so that a user can mimic interacting with the building. All of our data is formatted to be compliant with the popular open source Google Digital Buildings Ontology (DBO). The agent communicates with the building using the Protobuf open source serialization format(Google). The agent can send information requests to the building, asking for structural information, such as the number of devices, and telemetry information, such as the value of a particular sensor, and the building sends back a response, containing the requested information. The agent can also request that a setpoint be changed to a new value, and the building will respond if the change was successful.

Following the RL paradigm, the data in our dataset falls under the following categories:

1. **Environment Data** For each building environment, the dataset contains information on all HVAC zones and HVAC devices. For zones this includes the name and size of each zone, as well as how many devices are contained within it. For devices, this includes the zone the device is associated with, as well as every device sensor and setpoint.
2. **Observation Data** Observations consist of the measurements from all devices in the building (VAV’s zone air temperature, gas meter’s flow rate, etc.), provided at each time step.
3. **Action Data** The device setpoint values that the agent wants to set, provided at each timestep
4. **Reward Data** Information used to calculate the reward, as expressed in cost in dollars, carbon footprint, and comfort level of occupants, provided at each time step

The dataset currently consists of six years of data from three buildings. The details are in Table 1. For more details regarding the format of the data, including definitions and examples of each type of proto, see appendix B.

**Data Visualization** We also present a data visualization module, compatible both for viewing the real-world historical data, as well as visualizing the state of the simulator, as shown in Figure 1. Given an observation of a building environment, our visualization module renders a two dimensional heat-map view of the building. This greatly aids in understanding the data, and is invaluable in understanding how a particular policy is behaving.

Table 1: Building Information

BUILDING	FT <sup>2</sup>	FLOORS	DEVICES
SB1	93,858	2	170
SB2	62,613	1	152
SB3	118,086	3	152

## 5 SIMULATOR DESIGN CONSIDERATIONS

A fundamental trade-off when designing a simulator is speed versus fidelity, as depicted in Figure 2. Fidelity is the simulator’s ability to reproduce the building’s true dynamics that affect the optimization process. Speed refers to both simulator configuration time, i.e., the time required to configure a simulator for a target building, and the agent training time, i.e., the time necessary for the agent to optimize its policy using the simulator.

Every building is unique, due to its physical layout, equipment, and location. Fully customizing a high fidelity simulation to a specific target building requires nearly exhaustive knowledge of the building structure, materials, location, etc., some of which are unknowable, especially for legacy office buildings. This requires manual “guesstimation”, which can erode the accuracy promised by high-fidelity simulation. In general, the configuration time required for high-fidelity simulations limits their utility for deploying RL-based optimization to many buildings. High-fidelity simulations also are affected by computational demand and long execution times.

Alternatively, we propose a fast, low-to-medium-fidelity simulation model that was useful in addressing various design decisions, such as the reward function, the modeling of different algorithms. and for end-to-end testing. The simulation is built on a 2D finite-difference (FD) grid that models

thermal diffusion, and a simplified HVAC model that generates or removes heat on special “diffuser” control volumes (CV) in the FD grid. For more details on design considerations, see Appendix C.

While the uncalibrated simulator is of low-to-medium fidelity, the key additional factor is data. We collect recorded observations from the target building under baseline control, and use that data to **calibrate** the simulator, by adjusting the simulator’s physical parameters to minimize difference between real and simulated data. We believe this approach hits the sweet spot in this tradeoff, enabling scalability, while maintaining a high enough level of fidelity to train an improved policy.

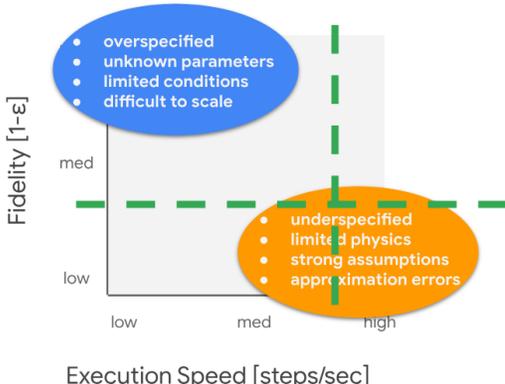


Figure 2: Simulation Fidelity vs. Execution Speed. The ideal operating point for training RL agents for energy and emission efficiency is a tradeoff between fidelity, depicted as 1 minus a normalized error  $\epsilon$  between simulation and real, and execution speed, as measured by the number of training steps per second. Additional consideration also includes the time to configure a custom simulator for the target building. While many approaches tend to favor high-fidelity over execution, speed, our approach argues a low-to-medium fidelity that has a medium-to-high speed is most suitable for training an RL agent.

## 6 A LIGHTWEIGHT, CALIBRATED SIMULATION

Our goal is to develop a method for applying RL at scale to commercial buildings. To this end, we put forth the following requirements for this to be feasible: We must have an easily customizable simulated environment to train the agent, with high enough fidelity to train an improved control agent. To meet these desiderata, we designed a light weight simulator based on finite differences approximation of heat exchange, building upon earlier work (Goldfeder & Sipple, 2023). We proposed a simple automated procedure to go from building floor plans to a custom simulator in a short time, and we designed a calibration and evaluation pipeline, to use data to fine tune the simulation to better match the real world. What follows is a description of our implementation.

**Thermal Model for the Simulation** As a template for developing simulators that represent target buildings, we start with a general-purpose high-level thermal model for simulating office buildings, illustrated in Figure 3. In this thermal cycle, we highlight significant energy consumers as follows. The boiler burns natural gas to heat the water,  $\dot{Q}_b$ . Water pumps consume electricity  $\dot{W}_{b,p}$  to circulate heating water through the VAVs. The air handler fans consume electricity  $\dot{W}_{b,in}$ ,  $\dot{W}_{b,out}$  to circulate the air through the VAVs. A motor drives the chiller’s compressor to operate a refrigeration cycle, consuming electricity  $\dot{W}_c$ . In some buildings coolant is circulated through the air handlers with pumps that consume electricity,  $\dot{W}_{c,p}$ .

We selected **water supply temperature**  $\hat{T}_b$  and the **air handler supply temperature**  $\hat{T}_s$  as agent actions because they affect the balance of electricity and natural gas consumption, they affect multiple device interactions, and they affect occupant comfort. Greater efficiencies can be achieved with these setpoints by choosing the ideal times and values to warm up and cool down the building in the workday mornings and evenings. Further tradeoffs include balancing the thermal load between hot water heating with natural gas and supply air heating with electricity using the air conditioner or heat pump units.

**Finite Differences Approximation** The diffusion of thermal energy in time and space of the building can be approximated using the method of Finite Differences (FD)(Sparrow, 1993; Lomax et al., 2002), and applying an energy balance. This method divides each floor of the building into a grid of three-dimensional control volumes and applies thermal diffusion equations to estimate the temperature of each control volume. By assuming each floor is adiabatically isolated, (i.e., no heat is transferred between floors), we can simplify the three-spatial dimensions into a spatial two-dimensional heat transfer problem. Each control volume is a narrow volume bounded horizontally, parameter-

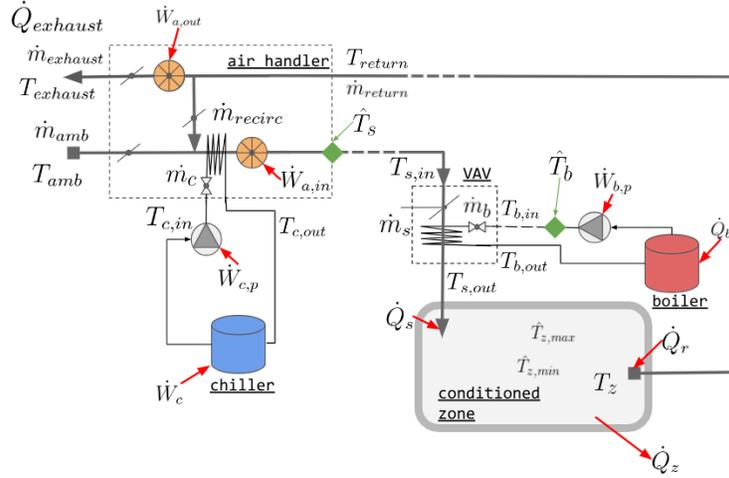


Figure 3: Thermal model for simulation. A building consists of conditioned zones, where the mean temperature of the zone  $T_z$  should be within upper and lower setpoints,  $\hat{T}_{z,max}$  and  $\hat{T}_{z,min}$ . Thermal power for heating or cooling the room is supplied to each zone,  $\dot{Q}_s$ , and recirculated from the zone,  $\dot{Q}_r$  from the HVAC system, with additional thermal exchange  $\dot{Q}_z$  from walls, doors, etc. The Air Handler supplies the building with air at supply air temperature setpoint  $\hat{T}_s$  drawing fresh air,  $\dot{m}_{amb}$ , at ambient temperatures,  $T_{amb}$ , and returning exhaust air  $\dot{m}_{exhaust}$  at temperature  $T_{exhaust}$  to the outside using intake and exhaust fans,  $\dot{W}_{a,in}$  and  $\dot{W}_{a,out}$ . A fraction of the return air can be recirculated,  $\dot{m}_{recirc}$ . Central air conditioning is achieved with a chiller and pump that joins a refrigeration cycle to the supply air, consuming electrical energy for the AC compressor  $\dot{W}_c$  and coolant circulation,  $\dot{W}_{c,p}$ . The hot water cycle consists of a boiler that maintains the supply water temperature at  $T_b$  heated by natural gas power  $\dot{Q}_b$ , and a pump that circulates hot water through the building, with electrical power  $\dot{W}_{b,p}$ . Supply air is delivered to the zones through Variable Air Volume (VAV) devices.

ized by  $\Delta x^2$ , and vertically by the height of the floor. The energy balance, shown below, is applied to each discrete control volume in the FD grid, and consists of the following components: (a) the thermal exchange across each face of the four participating faces control volume via conduction or convection  $Q_1, Q_2, Q_3, Q_4$ , (b) the change in internal energy over time in the control volume  $Mc \frac{\Delta T}{\Delta t}$ , and (c) an external energy source that enables applying local thermal energy from the HVAC model only for those control volumes that include an airflow diffuser,  $Q_{ext}$ . The equation is  $Q_{ext} + Q_1 + Q_2 + Q_3 + Q_4 = Mc \frac{\Delta T}{\Delta t}$ , where  $M$  is the mass and  $c$  is the heat capacity of the control volume,  $\Delta T$  is the temperature change from the prior timestep and  $\Delta t$  is the timestep interval.

The thermal exchange in (a) is calculated using Fourier’s law of steady conduction in the interior control volumes (walls and interior air), parameterized by the conductivity of the volume, and the exchange across the exterior faces of control volumes are calculated using the forced convection equation, parameterized by the convection coefficient, which approximates winds and currents surrounding the building. The change in internal energy (b) is parameterized by the density, and heat capacity of the control volume. Finally, the thermal energy associated with the VAV (c) is equally distributed to all associated control volumes that have a diffuser. Thermal diffusion within the building is mainly accomplished via forced or natural convection currents, which can be notoriously difficult to estimate accurately. We note that heat transfer using air circulation is effectively the exchange of air mass between control volumes, which we approximate by a randomized shuffling of air within thermal zones, parameterized by a shuffle probability and radius. For more details on this approximation and associated equations, see Appendix D.

**Simulator Configuration** For RL to scale to many buildings, it is critical to be able to easily and rapidly configure the simulator to any arbitrary building. We designed a procedure that, given floor-plans and HVAC layout information, enables generating a fully specified simulation very rapidly.

For example, on SB1, consisting of two floors and 170 devices, a single technician was able to configure the simulator in under three hours. Details of this procedure are provided in Appendix E.

**Simulator Calibration and Evaluation** In order to calibrate the simulator to the real world using data, we must have a metric with which to evaluate our simulator’s fidelity, and an optimization method to improve our simulator on this metric.

**$N$ -Step Evaluation** We propose a novel evaluation procedure, based on  $N$ -step prediction. Each iteration of our simulator was designed to represent a five-minute interval, and our real-world data is also obtained in five-minute intervals. To evaluate the simulator, we take a chunk of real data, consisting of  $N$  consecutive observations. We then initialize the simulator so that its initial state matches that of the starting observation, and run the simulator for  $N$  steps, replaying the same HVAC policy as was used in the real world. We then calculate our simulation fidelity metric, which is the mean absolute error of the temperatures in each temperature sensor at each time step, averaged over time. More formally, we define the Temporal Spatial Mean Absolute Error (TS-MAE) of  $Z$  zones over  $N$  timesteps as:

$$\epsilon = \sum_{t=1}^N \frac{1}{N} \left[ \frac{1}{Z} \sum_{z=1}^Z |T_{real,t,z} - T_{sim,t,z}| \right] \quad (1)$$

Where  $T_{real,t,z}$  is the measured zone air temperature for zone  $z$  at timestamp  $t$ , and  $T_{sim,t,z} = \frac{1}{|C_z|} \sum_{c=1}^{C_z} T_{t,c}$  is the mean temperature of all control volumes  $C_z$  in zone  $z$  at time  $t$ .

**Hyperparameter Calibration** Once we defined our simulation fidelity metric, the TS-MAE, we can attempt to minimize this error, thus improving fidelity, by hyperparameter tuning several physical constants and other variables using black-box optimization methods. We chose the method outlined in Golovin et. al. (Golovin et al., 2017), which automatically chooses the most appropriate strategy from a variety of popular algorithms.

## 7 SIMULATOR CALIBRATION

We now provide a full end-to-end demonstration of our calibration procedure, and show that our simulator, when tuned and calibrated, is able to make useful real-world predictions, and can train an RL agent to produce an improved policy over the baseline.

**Setup** We calibrated the simulator using data from SB1, with two stories, a combined surface area of 93,858 square feet, and 170 HVAC devices. Using the configuration pipeline, we went from floor plan blueprints to a fully configured simulator for this building, a process that took a single technician less than three hours to complete.

**Calibration Data** To calibrate our simulator, we took real-world data from three days, from Monday July 10, 2023 12:00 AM PST, to Thursday July 13, 2023 12:00 AM PST. The first two days were used as a train set, and the third day as validation of the calibrated performance on unseen data, as can be seen in Table 2. All times are given in US Pacific, the local time of the real building.

**Calibration Procedure** We ran hyperparameter tuning for 4000 iterations, with the aim of optimizing the TS-MAE, as outlined in equation 1, over the train data. We reviewed the physical constants that yielded the lowest simulation error from calibration. Densities, heat capacities, and conductivities plausibly matched common interior and exterior building materials. However, the external convection coefficient was higher than under the weather conditions, and likely is compensating for the radiative losses and gains, which were not directly simulated. For details about the hyperparameter tuning procedure, including the parameters varied, the ranges given, and the values found that best minimized the calibration metric, see Appendix F.

**Calibration Results** In Table 2, we present the predictive results of our calibrated simulator, on  $N$ -step prediction, for the train scenario, where  $N = 576$ , representing a two day predictive window, and the test scenario, where  $N = 288$ , representing a one day window. We calculated the TS-MAE, as defined in equation 1. We show results for the hyperparameters that best fit the train set, as well as for an uncalibrated simulator as a baseline. At no point was the validation data ever provided to

the tuning process. Note that the validation period is half the duration of the train period, so a lower error does not mean we are performing better than on the train data.

Table 2: Training and test data scenarios

SPLIT	LENGTH	START	END	CALIBRATED $\epsilon$	UNCALIBRATED $\epsilon$
TRAIN	48 HRS	2023-07-10 12AM	2023-07-12 12AM	0.717 $^{\circ}C$	1.971 $^{\circ}C$
VAL.	24 HRS	2023-07-12 12AM	2023-07-13 12AM	0.566 $^{\circ}C$	1.618 $^{\circ}C$

As indicated in Table 2, our tuning procedure drifts only 0.56  $^{\circ}C$  on average over a 24-hour period on the validation set.

**Visualizing Temperature Drift Over Time** Figure 4 illustrates temperature drift over time for the training scenario. At each time step, we calculate the spatial temperature for all sensors in both the real building and simulator, and present them as side-by-side boxplot distributions for comparison. Figure 5 shows the same for the validation scenario.

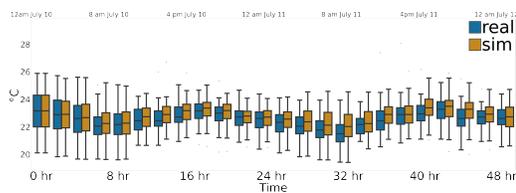


Figure 4: Drift Over 48 hrs on Train Set

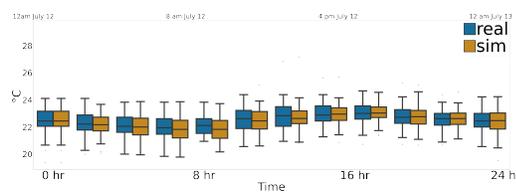


Figure 5: Drift Over 24 hrs on Validation Set

Here we can see that our simulator temperature distribution maintains a minimal drift from the real world, although it does seem a bit less reactive to daily fluctuation patterns, which may be the result of the lack of a radiative heat transfer model.

**Visualizing Spatial Errors** Figure 6 illustrates the results of this predictive process over a 24-hour period, on the validation data. It displays a heatmap of the spatial temperature difference throughout the building, between the real world and simulator, after 24 hours of the simulator making predictions. The ring of blue around the building indicates that our simulator is too cold on the perimeter, which implies that the heat exchange with the outside is happening more rapidly than it would in the real world. The inside of the building, at least on the first floor, contains significant amounts of red, indicating that despite the simulator perimeter being cooler than the real world, the inside is warmer. This implies that our thermal exchange within the building is not as rapid as that of the real world. We suspect that this may be because our simulator does not have a radiative heat transfer model. Lastly, there is a large amount of white in this image, indicating that for the most part, even after 24 hours of making predictions on the validation data, our calibration process was successful and the fidelity remains high. For more visuals of spatial errors, see appendix G.

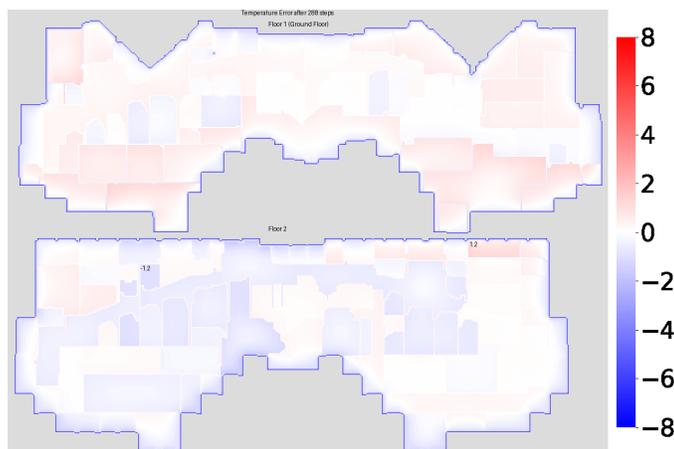


Figure 6: Visualization of simulator drift after 24 hours, on the validation data. The image is a heatmap representing the temperature difference between the simulator and the real world, with red indicating the simulator is hotter, blue indicating it is colder, and white indicating no difference. The zone with the max and min temperature difference are indicated by displaying above them the difference.

## 8 DEMONSTRATION BENCHMARKING RESULTS

While we believe our benchmark will be useful for a variety of tasks, such as further use of the data to calibrate the simulator, in this section we highlight results on three important tasks that our suite is well suited to: training an RL agent on the simulator, training a time-series regression model to predict the real world data, and training an RL agent on the real data directly.

**Training a Reinforcement Learning Agent on the Simulator** To demonstrate the usefulness of our calibrated simulator on generating an improved policy, we used Soft Actor Critic (SAC) algorithm (Haarnoja et al., 2018) to train an agent, and then compared our agent with the baseline performance of running the policy currently used in the real building. Both actor and critic were feedforward networks. We ran hyperparameter tuning, again using the method from Golovin et al. (Golovin et al., 2017), to choose the dimensionality of the critic network and actor network, the batch size, the critic learning rate and actor learning rate, and  $\gamma$ . We recorded the actor loss, critic loss, alpha loss, and return, over a two day period. The agents trained for 4,000 iterations. Using the  $R_{3C}$  reward, the baseline over this two day period had a return of -12.9, and our best agent had an improved return of -11.9, an 8% improvement over the baseline, as show in Table 3. For further training details, and an in depth performance comparison between the learned policy and the baseline, including a breakdown on setpoint deviation, carbon emissions, electrical energy, and natural gas energy, see Appendix H.

Table 3: Policy Comparison

POLICY	RETURN
BASELINE	-12.9
SAC	-11.9

**Training a Learned Dynamics Model** Another important task is to use a sequence model to learn to predict the real world data, effectively learning a dynamics model that can then be used in turn in place of the simulator to train an agent. To demonstrate this approach, we trained an encoder-decoder LSTM(Hochreiter, 1997) to model the building dynamics. The model takes in a historical sequence of length  $N$  and outputs a prediction sequence of length  $M$ . At each timestep  $t$  in the sequence, the model is given an observation  $O_t$ , action taken by the policy  $A_t$ , and auxiliary state features (such as time of day and weather, that are useful as inputs but need not be predicted)  $U_t$ , and for future timesteps, the model is trained to predict future observations, as well as future reward information (based on predicted energy use and carbon emissions)  $E_t$ . We evaluated this model by comparing its predictions with the real world data over a three week period, finding that it achieved strong performance and successfully modeled many building dynamics. For detailed architecture diagrams, training information and performance analysis, see Appendix I.

**Training a Reinforcement Learning Agent on Real Data** Building directly off of the above, we also trained an RL agent on the learned dynamics model, demonstrating the ability to learn a policy directly from data without involving the simulator. Like the simulator SAC agent, we were able to learn a policy that improved upon the baseline. For detailed analysis of this policy, see Appendix J.

## 9 LIMITATIONS AND CONCLUSION

The biggest limitation of our benchmark is that all buildings are located in California. We intend to remedy this in the near future by adding more buildings. Another limitation is that we only include data from a one year duration, and in the future we may add longer sequences, for year over year analysis. Our simulator also lacks a radiative heat model, and we hope further work can add this. In addition, our calibration focused on temperature, but in future work we hope to include energy consumption metrics as part of the calibration procedure.

We present a high quality interactive HVAC Control Suite, with real-world historical data from three buildings, as well as calibrated simulators for each building, and a novel, data-based, simulation calibration procedure. We also show promising initial results on key benchmark tasks. We believe this benchmark will facilitate collaboration, reproducibility, and progress on this problem, making an important contribution towards environmental sustainability.

## REFERENCES

Nilofar Asim, Marzieh Badiei, Masita Mohammad, Halim Razali, Armin Rajabi, Lim Chin Haw, and Mariyam Jameelah Ghazali. Sustainability of heating, ventilation and air-conditioning (hvac)

- 540 systems in buildings—an overview. *International journal of environmental research and public*  
541 *health*, 19(2):1016, 2022.
- 542
- 543 Donald Azuatalam, Wee-Lih Lee, Frits de Nijs, and Ariel Liebman. Reinforcement learning for  
544 whole-building hvac control and demand response. *Energy and AI*, 2:100020, 2020.
- 545 Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskiy,  
546 Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human benchmark.  
547 In *International conference on machine learning*, pp. 507–517. PMLR, 2020.
- 548
- 549 Craig Bakker, A August, S Vasisht, S Huang, and DL Vrabie. Data description for a high fidelity  
550 building emulator with building faults. *tech. rep.*, 2022.
- 551 Yasaman Balali, Adrian Chong, Andrew Busch, and Steven O’Keefe. Energy modelling and control  
552 of building heating and cooling systems with data-driven and hybrid models—a review. *Renew-*  
553 *able and Sustainable Energy Reviews*, 183:113496, 2023.
- 554
- 555 Sean Barker, Aditya Mishra, David Irwin, Emmanuel Cecchet, Prashant Shenoy, Jeannie Albrecht,  
556 et al. Smart\*: An open data set and tools for enabling research in sustainable homes. *SustKDD*,  
557 *August*, 111(112):108, 2012.
- 558 Mangesh Basarkar. Modeling and simulation of hvac faults in energyplus. 2011.
- 559
- 560 Marco Biemann, Fabian Scheller, Xiufeng Liu, and Lizhen Huang. Experimental evaluation of  
561 model-free reinforcement learning algorithms for continuous hvac control. *Applied Energy*, 298:  
562 117164, 2021.
- 563 Christian Blad, Simon Bøgh, and Carsten Skovmose Kallesøe. Data-driven offline reinforcement  
564 learning for hvac-systems. *Energy*, 261:125290, 2022.
- 565
- 566 Michael Blonsky, Jeff Maguire, Killian McKenna, Dylan Cutler, Sivasathya Pradha Balamurugan,  
567 and Xin Jin. Ochre: The object-oriented, controllable, high-resolution residential energy model  
568 for dynamic integration studies. *Applied Energy*, 290:116732, 2021.
- 569 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
570 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 571
- 572 Gustavo Campos, Nael H El-Farra, and Ahmet Palazoglu. Soft actor-critic deep reinforcement  
573 learning with hybrid mixed-integer actions for demand responsive scheduling of energy systems.  
574 *Industrial & Engineering Chemistry Research*, 61(24):8443–8461, 2022.
- 575 Bingqing Chen, Zicheng Cai, and Mario Bergés. Gnu-rl: A practical and scalable reinforcement  
576 learning solution for building hvac control using a differentiable mpc policy. *Frontiers in Built*  
577 *Environment*, 6:562239, 2020.
- 578 Liangliang Chen, Fei Meng, and Ying Zhang. Fast human-in-the-loop control for hvac systems via  
579 meta-learning and model-based offline reinforcement learning. *IEEE Transactions on Sustainable*  
580 *Computing*, 2023.
- 581
- 582 Davide Coraci, Silvio Brandi, Marco Savino Piscitelli, and Alfonso Capozzoli. Online implemen-  
583 tation of a soft actor-critic agent to enhance indoor temperature control and energy efficiency in  
584 buildings. *Energies*, 14(4):997, 2021.
- 585 Drury B Crawley, Linda K Lawrie, Frederick C Winkelmann, Walter F Buhl, Y Joe Huang, Curtis O  
586 Pedersen, Richard K Strand, Richard J Liesen, Daniel E Fisher, Michael J Witte, et al. Energyplus:  
587 creating a new-generation building energy simulation program. *Energy and buildings*, 33(4):319–  
588 331, 2001.
- 589
- 590 Ján Drgoňa, Aaron R Tuor, Vikas Chandan, and Draguna L Vrabie. Physics-constrained deep learn-  
591 ing of multi-zone building thermal dynamics. *Energy and Buildings*, 243:110992, 2021.
- 592 EIA. Frequently Asked Questions (FAQs) - U.S. Energy Information Administration (EIA) —  
593 eia.gov. <https://www.eia.gov/tools/faqs/faq.php?id=86&t=1>. [Accessed 04-06-2024].

- 594 Xi Fang, Guangcai Gong, Guannan Li, Liang Chun, Pei Peng, Wenqiang Li, Xing Shi, and Xiang  
595 Chen. Deep reinforcement learning optimal control strategy for temperature setpoint real-time  
596 reset in multi-zone building hvac system. *Applied Thermal Engineering*, 212:118552, 2022.  
597
- 598 Kristin Field, Michael Deru, and Daniel Studer. Using doe commercial reference buildings for  
599 simulation studies. 2010.
- 600 Kwong Fai Fong, Victor Ian Hanby, and Tin-Tai Chow. Hvac system optimization for energy man-  
601 agement by evolutionary programming. *Energy and buildings*, 38(3):220–231, 2006.  
602
- 603 Cheng Gao and Dan Wang. Comparative study of model-based and model-free reinforcement learn-  
604 ing control performance in hvac systems. *Journal of Building Engineering*, 74:106852, 2023.
- 605 Frédéric Garcia and Emmanuel Rachelson. Markov decision processes. *Markov Decision Pro-*  
606 *cesses in Artificial Intelligence*, pp. 1–38, 2013.  
607
- 608 Judah A Goldfeder and John A Sipple. A lightweight calibrated simulation enabling efficient offline  
609 learning for optimal control of real buildings. In *Proceedings of the 10th ACM International*  
610 *Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*, pp. 352–356,  
611 2023.
- 612 Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and David Scul-  
613 ley. Google vizier: A service for black-box optimization. In *Proceedings of the 23rd ACM*  
614 *SIGKDD international conference on knowledge discovery and data mining*, pp. 1487–1495,  
615 2017.
- 616 Google. Protocol buffers. <http://code.google.com/apis/protocolbuffers/>.  
617
- 618 Jessica Granderson, Guanqing Lin, Yimin Chen, Armando Casillas, Jin Wen, Zhelun Chen, Piljae  
619 Im, Sen Huang, and Jiazhen Ling. A labeled dataset for building hvac systems operating in faulted  
620 and fault-free states. *Scientific data*, 10(1):342, 2023.
- 621 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash  
622 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-  
623 cations. *arXiv preprint arXiv:1812.05905*, 2018.  
624
- 625 Philipp Heer, Curdin Derungs, Benjamin Huber, Felix Bünning, Reto Fricker, Sascha Stoller, and  
626 Björn Niesen. Comprehensive energy demand and usage data for building automation. *Scientific*  
627 *Data*, 11(1):469, 2024.
- 628 S Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.  
629
- 630 A Husaunndee, R Lahrech, H Vaezi-Nejad, and JC Visier. Simbad: A simulation toolbox for the de-  
631 sign and test of hvac control systems. In *Proceedings of the 5th international IBPSA conference*,  
632 volume 2, pp. 269–276. International Building Performance Simulation Association (IBPSA)  
633 Prague . . . , 1997.
- 634 Farrokh Jazizadeh, Milad Afzalan, Burcin Becerik-Gerber, and Lucio Soibelman. Embed: A dataset  
635 for energy monitoring through building electricity disaggregation. In *Proceedings of the Ninth*  
636 *International Conference on Future Energy Systems*, pp. 230–235, 2018.  
637
- 638 Anjukan Kathirgamanathan, Eleni Mangina, and Donal P Finn. Development of a soft actor critic  
639 deep reinforcement learning approach for harnessing energy flexibility in a large office building.  
640 *Energy and AI*, 5:100101, 2021.
- 641 Dongsu Kim, Jongman Lee, Sunglok Do, Pedro J Mago, Kwang Ho Lee, and Heejin Cho. Energy  
642 modeling and model predictive control for hvac in buildings: a review of current research trends.  
643 *Energies*, 15(19):7231, 2022.
- 644 Peter Klanatsky, François Veynandt, and Christian Heschl. Grey-box model for model predictive  
645 control of buildings. *Energy and Buildings*, 300:113624, 2023.  
646
- 647 Thomas Kriechbaumer and Hans-Arno Jacobsen. Blond, a building-level office environment dataset  
of typical electrical appliances. *Scientific data*, 5(1):1–14, 2018.

- 648 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 649
- 650
- 651 Harvard Lomax, Thomas H Pulliam, David W Zingg, and TA Kowalewski. Fundamentals of computational fluid dynamics. *Appl. Mech. Rev.*, 55(4):B61–B61, 2002.
- 652
- 653 Na Luo, Zhe Wang, David Blum, Christopher Weyandt, Norman Bourassa, Mary Ann Piette, and Tianzhen Hong. A three-year dataset supporting research on building energy management and occupancy analytics. *Scientific data*, 9(1):156, 2022.
- 654
- 655
- 656
- 657 Karl Mason and Santiago Grijalva. A review of reinforcement learning for autonomous building energy management. *Computers & Electrical Engineering*, 78:300–312, 2019.
- 658
- 659 Paul A Mathew, Laurel N Dunn, Michael D Sohn, Andrea Mercado, Claudine Custudio, and Travis Walter. Big-data for building energy performance: Lessons from assembling a very large national database of building energy use. *Applied Energy*, 140:85–93, 2015.
- 660
- 661
- 662
- 663 Faye C McQuiston, Jerald D Parker, Jeffrey D Spitler, and Hessam Taherian. *Heating, ventilating, and air conditioning: analysis and design*. John Wiley & Sons, 2023.
- 664
- 665 Christoph J Meinrenken, Noah Rauschkolb, Sanjmeet Abrol, Tuhin Chakrabarty, Victor C Decalf, Christopher Hidey, Kathleen McKeown, Ali Mehmani, Vijay Modi, and Patricia J Culligan. Mfred, 10 second interval real and reactive power for groups of 390 us apartments of varying size and vintage. *Scientific Data*, 7(1):375, 2020.
- 666
- 667
- 668
- 669
- 670 Clayton Miller, Anjukan Kathirgamanathan, Bianca Picchetti, Pandarasamy Arjunan, June Young Park, Zoltan Nagy, Paul Raftery, Brodie W Hobson, Zixiao Shi, and Forrest Meggers. The building data genome project 2, energy meter data from the ashrae great energy predictor iii competition. *Scientific data*, 7(1):368, 2020.
- 671
- 672
- 673
- 674 Michael C Mozer. The neural network house: An environment that adapts to its inhabitants. In *Proc. AAAI Spring Symp. Intelligent Environments*, volume 58, pp. 110–114, 1998.
- 675
- 676
- 677 David Murray, Lina Stankovic, and Vladimir Stankovic. An electrical load measurements dataset of united kingdom households from a two-year longitudinal study. *Scientific data*, 4(1):1–12, 2017.
- 678
- 679 Biswas Gautam Naug, Avisek and Vikas. Chandan. Vanderbilt alumni hall, nashville, tennessee. URL <https://bbd.labworks.org/ds/vah>.
- 680
- 681
- 682 Cheol Park, Daniel R Clark, and George E Kelly. An overview of hvacsim+, a dynamic building/hvac/control systems simulation program. In *Proceedings of the 1st Annual Building Energy Simulation Conference, Seattle, WA*, pp. 21–22, 1985.
- 683
- 684
- 685 Luis Pérez-Lombard, José Ortiz, and Christine Pout. A review on buildings energy consumption information. *Energy and buildings*, 40(3):394–398, 2008.
- 686
- 687
- 688 Betsy Pettit, Cathy Gates, A Hunter Fanney, and William Healy. Design challenges of the net zero energy residential test facility. *Gaithersburg, MD: National Institute of Standards and Technology*, 2014.
- 689
- 690
- 691 H Rashid, P Singh, and A Singh. I-blend, a campus-scale commercial and residential buildings electrical energy dataset. *scientific data*, 6 (1), 2019.
- 692
- 693
- 694 Peter Riederer. Matlab/simulink for building and hvac simulation-state of the art. In *Ninth International IBPSA Conference*, pp. 1019–1026, 2005.
- 695
- 696 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 697
- 698
- 699
- 700 Olga Sachs, Verena Tiefenbeck, Caroline Duvier, Angela Qin, Kate Cheney, Craig Akers, and Kurt Roth. Field evaluation of programmable thermostats. Technical report, National Renewable Energy Lab.(NREL), Golden, CO (United States), 2012.
- 701

- 702 Igor Sartori, Harald Taxt Walnum, Kristian S Skeie, Laurent Georges, Michael D Knudsen, Peder  
703 Bacher, José Candanedo, Anna-Maria Sigounis, Anand Krishnan Prakash, Marco Pritoni, et al.  
704 Sub-hourly measurement datasets from 6 real buildings: Energy use and indoor climate. *Data in*  
705 *Brief*, 48:109149, 2023.
- 706 R Sendra-Arranz and A Gutiérrez. A long short-term memory artificial neural network to predict  
707 daily hvac consumption in buildings. *Energy and Buildings*, 216:109952, 2020.
- 708 Olli Seppanen, William J Fisk, and QH Lei. Effect of temperature on task performance in office  
709 environment. 2006.
- 710 E M Sparrow. Heat transfer: Conduction [lecture notes], 1993.
- 711 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 712 Saman Taheri, Paniz Hosseini, and Ali Razban. Model predictive control of heating, ventilation, and  
713 air conditioning (hvac) systems: A state-of-the-art review. *Journal of Building Engineering*, 60:  
714 105067, 2022.
- 715 TFDS. TensorFlow Datasets, a collection of ready-to-use datasets. [https://www.](https://www.tensorflow.org/datasets)  
716 [tensorflow.org/datasets](https://www.tensorflow.org/datasets).
- 717 Marija Trčka and Jan LM Hensen. Overview of hvac system simulation. *Automation in construction*,  
718 19(2):93–99, 2010.
- 719 Marija Trcka, Michael Wetter, and JLM Hensen. Comparison of co-simulation approaches for build-  
720 ing and hvac/r system simulation. In *10th International IBPSA Building Simulation Conference*  
721 *(BS 2007), September 3-6, 2007, Beijing, China*, pp. 1418–1425, 2007.
- 722 Marija Trčka, Jan LM Hensen, and Michael Wetter. Co-simulation of innovative integrated hvac  
723 systems in buildings. *Journal of Building Performance Simulation*, 2(3):209–230, 2009.
- 724 Bryan Urban, Kurt Roth, Olga Sachs, Verena Tiefenbeck, Caroline Duvier, Kate Cheney, and Craig  
725 Akers. Multifamily programmable thermostat data. doi: 10.25984/1844177.
- 726 José R Vázquez-Canteli and Zoltán Nagy. Reinforcement learning for demand response: A review  
727 of algorithms and modeling techniques. *Applied energy*, 235:1072–1089, 2019.
- 728 Kirubakaran Velswamy, Biao Huang, et al. A long-short term memory recurrent neural network  
729 based reinforcement learning controller for office heating ventilation and air conditioning systems.  
730 2017.
- 731 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
732 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*  
733 *preprint arXiv:1804.07461*, 2018.
- 734 Marshall Wang, John Willes, Thomas Jiralerspong, and Matin Moezzi. A comparison of classical  
735 and deep reinforcement learning methods for hvac control. *arXiv preprint arXiv:2308.05711*,  
736 2023.
- 737 Mubashir Wani, Akshya Swain, and Abhisek Ukil. Control strategies for energy optimization of  
738 hvac systems in small office buildings using energypplus tm. In *2019 IEEE Innovative Smart Grid*  
739 *Technologies-Asia (ISGT Asia)*, pp. 2698–2703. IEEE, 2019.
- 740 Tianshu Wei, Yanzhi Wang, and Qi Zhu. Deep reinforcement learning for building hvac control. In  
741 *Proceedings of the 54th annual design automation conference 2017*, pp. 1–6, 2017.
- 742 Yunyang Ye, Wangda Zuo, and Gang Wang. A comprehensive review of energy-related data for us  
743 commercial buildings. *Energy and Buildings*, 186:126–137, 2019.
- 744 Liang Yu, Yi Sun, Zhanbo Xu, Chao Shen, Dong Yue, Tao Jiang, and Xiaohong Guan. Multi-agent  
745 deep reinforcement learning for hvac control in commercial buildings. *IEEE Transactions on*  
746 *Smart Grid*, 12(1):407–419, 2020.

- 756 Liang Yu, Shuqi Qin, Meng Zhang, Chao Shen, Tao Jiang, and Xiaohong Guan. A review of deep  
757 reinforcement learning for smart building energy management. *IEEE Internet of Things Journal*,  
758 8(15):12046–12063, 2021.
- 759 Chi Zhang, Sanmukh R Kuppannagari, Rajgopal Kannan, and Viktor K Prasanna. Building hvac  
760 scheduling using reinforcement learning via neural network based model approximation. In *Pro-  
761 ceedings of the 6th ACM international conference on systems for energy-efficient buildings, cities,  
762 and transportation*, pp. 287–296, 2019a.
- 764 Chi Zhang, Yuanyuan Shi, and Yize Chen. Bear: Physics-principled building environment for con-  
765 trol and reinforcement learning. In *Proceedings of the 14th ACM International Conference on  
766 Future Energy Systems*, pp. 66–71, 2023.
- 767 Zhiang Zhang, Adrian Chong, Yuqi Pan, Chenlu Zhang, and Khee Poh Lam. Whole building energy  
768 model for hvac optimal control: A practical framework based on deep reinforcement learning.  
769 *Energy and Buildings*, 199:472–490, 2019b.
- 771 Huan Zhao, Junhua Zhao, Ting Shu, and Zibin Pan. Hybrid-model-based deep reinforcement learn-  
772 ing for heating, ventilation, and air-conditioning control. *Frontiers in Energy Research*, 8:610518,  
773 2021.
- 774 Jie Zhao, Khee Poh Lam, B Erik Ydstie, and Omer T Karaguzel. Energyplus model-based predic-  
775 tive control within design–build–operate energy information modelling infrastructure. *Journal of  
776 Building Performance Simulation*, 8(3):121–134, 2015.
- 777 Dian Zhuang, Vincent JL Gan, Zeynep Duygu Tekler, Adrian Chong, Shuai Tian, and Xing Shi.  
778 Data-driven predictive control for smart hvac system in iot-integrated buildings with time-series  
779 forecasting and reinforcement learning. *Applied Energy*, 338:120936, 2023.
- 781 Zhengbo Zou, Xinran Yu, and Semiha Ergan. Towards optimal control of air handling units us-  
782 ing deep reinforcement learning and recurrent neural network. *Building and Environment*, 168:  
783 106535, 2020.

## 785 A REWARD FUNCTION DETAILS

787 We call our reward function the 3C Reward, because it is made up of a combination of three fac-  
788 tors: Comfort, Cost, and Carbon. The purpose of the reward function is to provide the agent a  
789 feedback signal after each action about the quality of the current and past actions performed. We  
790 combine the different objectives described in Optimization Problem as a normalized, weighted sum  
791 of maintaining comfort conditions, electrical cost, and carbon cost:

$$792 R_{3C} = u \times C_1 + v \times C_2 + w \times C_3$$

795 where  $C_1$  represents normalized comfort conditions,  $C_2$  normalized energy cost and  $C_3$  normalized  
796 carbon emission. Constants  $u$ ,  $v$ ,  $w$  represent operator preferences, allowing them to weight the  
797 relative importance of cost, comfort and carbon consumption.

798 Each value  $C_1, C_2, C_3$ , is bounded by the range  $[-1, 0]$ , where worst performance is  $-1$  and the  
799 ideal performance upper-bound is 0 Thus the reward function in an aggregate is formulated as an  
800 approximate regret function, bounded in the range  $[-1, 0]$ , and represents an offset from the best-case  
801 where comfort conditions are perfectly maintained, without consuming energy and emitting carbon.  
802 Each of the sub functions  $C_1, C_2, C_3$  will be elaborated next.

### 803 A.1 COMFORT LOSS FUNCTION ( $C_1$ )

804 Besides zone air temperature, other factors such as ventilation, drafts, solar exposure, humidity  
805 and air quality affect human comfort and productivity in office buildings. However, for now we  
806 are focused solely on temperature as the indicator of the comfort level in the office buildings. As  
807 additional sensors are deployed and the other factors are measured, they should be considered in the  
808 definition of an enhanced comfort loss function.  
809

810 Studies have shown that a relationship exists between work performance and temperature. For ex-  
 811 ample, in Seppänen, et al. 2006 (Seppanen et al., 2006), work performance was quantified as the  
 812 mean time required to complete common office tasks (e.g., text processing, bookkeeping calcula-  
 813 tions, telephone customer service calls, etc.). Performance was shown to increase gradually with  
 814 temperatures increasing up to 21-22°C and decreasing at temperatures beyond 23-24°C. Therefore,  
 815 when temperatures deviate outside setpoints, the comfort loss should also be smooth and monoton-  
 816 ically increasing.

817 Thus, the following rules were selected to govern the comfort loss function:

- 818 1. Setpoints define the comfort standards, and no penalty should be applied whenever the zone  
 819 temperature is within heating and cooling setpoints.
- 820 2. Comfort is undefined when the zone is unoccupied: if the zone is unoccupied, comfort loss  
 821 is zero, regardless of zone temperature.
- 822 3. Comfort decays smoothly and monotonically as the temperatures drift from setpoints, and  
 823 occupants are tolerant to small setpoint deviations. Therefore, small setpoint deviations  
 824 should have a small comfort penalty, and the penalty should smoothly increase as the devi-  
 825 ations increase.
- 826 4. Large setpoint deviations should approach a maximum, bounded penalty, where a zone  
 827 becomes completely intolerable for its occupants.

828 The comfort loss function represents a bounded penalty term for occupied zones that have zone air  
 829 temperatures outside of setpoint, and covers three adjacent temperature intervals: below cooling  
 830 setpoint  $T_z < \hat{T}_{heating}$ , inside setpoints  $\hat{T}_{heating} \leq T_z \leq \hat{T}_{cooling}$ , and above cooling setpoint  
 831  $\hat{T}_{cooling} < T_z$

832 We propose a logistic sigmoid parameterized by  $\lambda$  and  $\Delta$  to represent the smooth decay (increase  
 833 loss) of comfort below the heating and above the cooling setpoints. Parameter  $\lambda$  is a stiffness coef-  
 834 ficient that affects the slope of the decay and parameter  $\Delta$  represents the offset in °C from the set  
 835 point where halfway loss value (0.5) occurs. Additionally we define a step function  $\delta(k) = 1$  when  
 836 the zone has at least one occupant ( $k > 0$ ), and  $\delta(k) = 0$  otherwise.

$$837 h_z(T_z, k_z, \hat{T}_{heating}, \hat{T}_{cooling}) = \begin{cases} \frac{\delta(k_z)}{1+e^{-\lambda(T_z-\hat{T}_{heating}+\Delta)}} - 1 & T_z < \hat{T}_{heating} \\ 0 & \hat{T}_{heating} \leq T_z \leq \hat{T}_{cooling} \\ \frac{-\delta(k_z)}{1+e^{-\lambda(T_z-\hat{T}_{cooling}-\Delta)}} & \hat{T}_{cooling} < T_z \end{cases}$$

838 The chart below shows the comfort loss curve with common setpoints, where the horizontal axis  
 839 represents zone air temperature and the vertical axis represents the loss. The heating and cooling  
 840 setpoints were taken from data recordings.

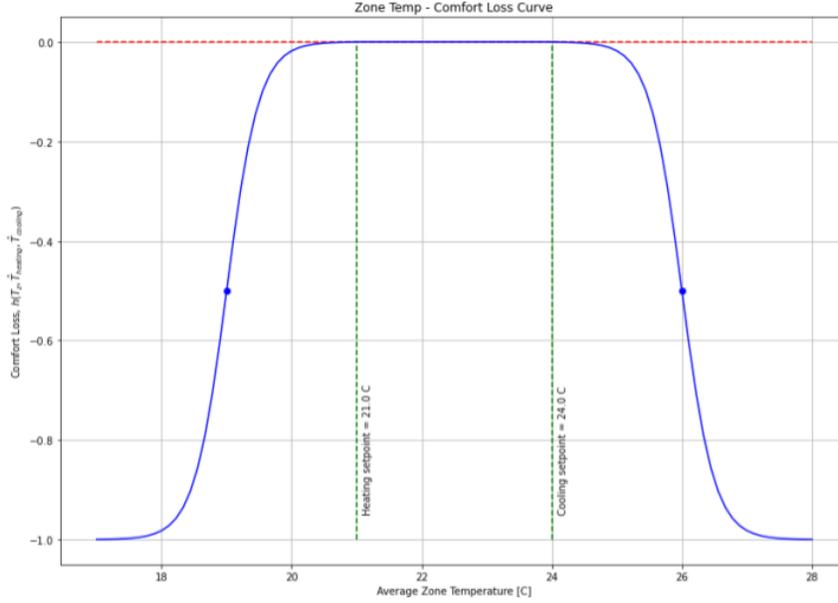


Figure 7: Setpoint Diagram

Finally, we compute the average of all zone comfort losses as the building’s overall comfort loss:

$$h_t(S_t) = \frac{1}{|Z|} \times \sum_{z \in Z} h_z(T_z, k, \hat{T}_{heating}, \hat{T}_{cooling})$$

**Live Occupant Feedback** The idea of human feedback shaping the agent’s policy may be particularly suitable for the smart buildings project, and has been detailed in Knox and Stone 2009. While not implemented in the initial version of the reward function, the comfort loss function can be extended with an occupant feedback signal reflecting discomfort (e.g., “too hot” or “too cold”) in a variety of methods like Mozer 1998 (Mozer, 1998). The agent’s goal should be to minimize this type of feedback, and the regret should be increased anytime this feedback signal is received. Suppose one or more occupants in zone  $z$ , provided a “too cold” feedback signal,  $\hat{T}_{heating}$  may be increased by a small amount from the baseline setpoint configuration, and may smoothly return to the baseline smoothly after an appropriate delay.

**Stochastic Occupancy Model** The occupancy signal  $k_z$  is the average number of occupants in zone  $z$  during a time step  $t_i - t_{i-1}$  and is used in computing the comfort loss function described above. Ideally, the occupancy signal is obtained from motion detection sensors or secondary indicators of occupancy, such as wifi signals, badge swipes, calendar appointments, etc. However, a data-driven occupancy signal was not available for the initial dataset, and the following stochastic occupancy model is used instead.

For workdays, we would like model occupancy as a process in the zone where a max number of occupants,  $k_{z,max}$  arrive at random times in an arrival window  $[\tau_{in,start}, \tau_{in,end}]$ , and depart the zone in a departure window  $[\tau_{out,start}, \tau_{out,end}]$ . The arrivals and departures should occur evenly within the intervals and the expectation of the arrival time should be at the halfway point of the arrival interval:

$$\mathbb{E}[\text{occupant arrival time}] = \frac{1}{2}(\tau_{in,end} - \tau_{in,start}) + \tau_{in,start}$$

Likewise, the expectation of the departure time should be at the halfway point of the departure interval:

$$\mathbb{E}[\text{occupant departure time}] = \frac{1}{2}(\tau_{out,end} - \tau_{out,start}) + \tau_{out,start}$$

If the number of timesteps within the arrival and departure intervals is  $n_{arrival}$  and  $n_{departure}$ , this process can be modeled as a geometric distribution where each timestep and occupant is a Bernoulli trial with probabilities:

$P(\text{occupant arrives} \mid \text{occupant has not yet arrived}) = \frac{2}{n_{arrival}}$  and  $P(\text{occupant departs} \mid \text{occupant has arrived}) = \frac{2}{n_{departure}}$ . During holidays and weekends, the zones are not occupied:  $k_z = 0$ .

## A.2 ENERGY COST FUNCTION ( $C_2$ )

The energy cost function  $C_1(S_t)$  is a normalized, aggregate cost estimate from consuming electrical and natural gas energy during one timestep. The cost function is the ratio of the actual energy used to the maximum energy capacity that ranges between 0: no cost incurred; and 1: maximum cost incurred.

$$C_2(S_t) = \frac{\text{actual energy cost}}{\text{cost at max energy capacity}}$$

General energy cost can be calculated as the product of the mean power applied, the time interval, and the cost per unit energy at the time of the interval, where we use  $W, \dot{W}$  to represent electrical/mechanical energy, and power, and  $Q, \dot{Q}$  to represent thermal energy and power from natural gas. Since all four terms contain the same interval  $t_i - t_{i-1}$ , they cancel out, allowing us to use power instead of energy. As described above, pumps, blowers, and AC/refrigeration cycles consume electricity and water heaters/boilers consume natural gas. Therefore the total energy and cost is the sum of each energy consumer cost used over the interval:

$$C_2(S_t) = \frac{(\dot{W}_a + \dot{W}_m + \dot{W}_p) \times p_e(t) + \dot{Q}_g \times p_g(t)}{(\dot{W}_{a,max} + \dot{W}_{m,max} + \dot{W}_{p,max}) \times p_e(t) + \dot{Q}_{g,max} \times p_g(t)}$$

Where  $\dot{W}_a$  and  $\dot{W}_{a,max}$  are the actual and max electrical power for the AC/refrigeration cycle,  $\dot{W}_m$  and  $\dot{W}_{m,max}$  are the actual and max electrical power for the blowers/air circulation,  $\dot{W}_p$  and  $\dot{W}_{p,max}$  are the actual and max pump electrical power, and  $\dot{Q}_g$  and  $\dot{Q}_{g,max}$  are the actual and max thermal power. Terms  $p_e(t)$  and  $p_g(t)$  are the electricity and gas price per energy incurred over the interval at time  $t$ .

The actual power terms in the numerator are estimated from the device observations and the device's fixed parameters using standard HVAC energy conversions. The max power terms in the denominator are derived from device ratings, which define the maximum operating nouns of the device.

## A.3 CARBON EMISSION COST FUNCTION ( $C_3$ )

Similar to the energy cost function, carbon emission cost function is a function of the electrical and natural gas power used during the interval. The carbon emission cost function  $C_3$  is a normalized, aggregate cost estimate from the emission of carbon mass by consuming electrical and natural gas energy during one timestep. The cost function is the ratio of the actual carbon used to the maximum carbon emitted that ranges between 0: no emission cost incurred; and 1: maximum emission cost incurred.

$$C_3(S_t) = \frac{\text{actual carbon mass emitted}}{\text{maximum carbon emitted}}$$

The carbon emission cost is similar to the energy cost function described above, except that we replace the price terms  $p_e, p_g$  with emission terms  $r_e, r_g$  that convert the power to carbon emission rates.

972 While the emission rate for natural gas is fairly constant, the emission rate for electricity is dependent  
 973 on the utility’s current renewable energy supply and consumer load during the interval and may  
 974 fluctuate significantly.  
 975

#### 976 A.4 IMMEDIATE AND DELAYED REWARD RESPONSES 977

978 The reward function is a weighted average of maintaining temperature setpoints in occupied zones,  
 979 while minimizing energy cost, and minimizing carbon emission. Both energy and carbon emission  
 980 cost functions provide a low latency response, because actions have an almost immediate effect on  
 981 the reward. For example, lowering the supply water temperature setpoint will reduce the flow of  
 982 natural gas to the burner, bringing  $\dot{Q}$  down in the next step. However, the effect of increasing water  
 983 temperature on the comfort loss function may be delayed by multiple time steps, due to the thermal  
 984 latency in the building. This thermal latency is due to inherent heat capacity and thermal resistance  
 985 within the building that has a dampening effect on diffusing heat throughout the building. This  
 986 means that some settings of  $u, v, w$  may cause undesirable effects. Experiments with the simulation  
 987 indicate that too strong weights (e.g.,  $u + v \geq 0.6$ ) toward energy cost and/or carbon emission may  
 988 lead the agent to lower the water temperature, which can cause the VAVs to increase their airflow  
 989 demand to compensate for a lower supply air temperature, since thermal energy flow is a tradeoff  
 990 between air mass flow and water heating at the VAV’s heat exchanger. Consequently, the increased  
 991 airflow demand results in a much higher, delayed electrical energy consumption by the blowers to  
 992 meet the zone airflow demand.  
 993

## 994 B PROTO DEFINITIONS 995

996 Here, we will elaborate on the exact proto definitions used in the dataset.

997 Having applied the RL paradigm, the data in our dataset falls under the following categories:  
 998

- 999 1. **Environment Data** General information about the environment, such as the number of  
 1000 devices and zones, and their names and device types. This is provided once per building  
 1001 environment
- 1002 2. **Observation Data** The measurements from all devices in the building (VAV’s zone air  
 1003 temperature, gas meter’s flow rate, etc.), provided at each time step
- 1004 3. **Action Data** The device setpoint values that the agent wants to set, provided at each  
 1005 timestep
- 1006 4. **Reward Data** Information used to calculate the reward, as expressed in energy cost in  
 1007 dollars, carbon emission, and comfort level of occupants, provided at each time step  
 1008  
 1009

1010 As mentioned above, this data is stored in protos. This section provides the definition of each proto,  
 1011 categorizing them using the four categories above, with examples of each.  
 1012

### 1013 B.1 ENVIRONMENT DATA PROTOS 1014

1015 This is the data that provides, once per environment, details about the environment such as number  
 1016 of devices, and zones, etc. There are two proto definitions:  
 1017

- 1018 1. **ZoneInfo:** The `ZoneInfo` message defines thermal spaces or zones in the building and  
 1019 provides zone-to-device association, which enables using the associated VAVs’ zone air  
 1020 temperatures to estimate the zone’s temperature.
- 1021 2. **DeviceInfo:** The HVAC devices in the building are defined in the `DeviceInfo` mes-  
 1022 sage. Each device exposes a map of `observable_fields` and `action_fields`. The  
 1023 `observable_fields` represent the observable state of the building in native units, and  
 1024 the `action_fields` are available setpoints exposed by the building that the agent may  
 1025 add to its action space. Currently `observable_fields` and `action_fields` are  
 floating point values, but may be expanded to categorical values in the future.

### 1026 B.1.1 ZONEINFO DEFINITION

```

1027
1028 1 message ZoneInfo {
1029 2
1030 3
1031 4 enum ZoneType {
1032 5     UNDEFINED = 0;
1033 6     ROOM = 1;
1034 7     FLOOR = 2;
1035 8     OTHER = 10;
1036 9 }
1037 10 // Unique Identifier of the zone.
1038 11 string zone_id = 1;
1039 12 // ID of the building
1040 13 string building_id = 2;
1041 14 // Free-form description of the zone, like microkitchen, office, etc.
1042 15 string zone_description = 3;
1043 16 // Square footage of the zone.
1044 17 float area = 4; // square meters
1045 18 // Zero to multiple device identifiers associated with this zone, like←
1046 19     VAVs.
1047 20 repeated string devices = 5;
1048 21 // Optional field to describe the type of zone.
1049 22 ZoneType zone_type = 6;
1050 23 // Optional field to indicate the floor of the building.
1051 24 int32 floor = 7;
1052 }

```

### 1050 B.1.2 ZONEINFO EXAMPLE

```

1052 1 zone_id: "rooms/9028552253"
1053 2 building_id: "buildings/3616672508"
1054 3 zone_description: "US-BLDG-2-C201"
1055 4 devices: "2614466029028994"
1056 5 devices: "2687242320524339"
1057 6 devices: "2640423556868160"
1058 7 zone_type: ROOM
1059 8 floor: 2

```

### 1060 B.1.3 DEVICEINFO DEFINITION

```

1062 1 // Details about a specific device in the building.
1063 2 message DeviceInfo {
1064 3 // Device types in smart buildings (official Carson top-level device ←
1065 4     types).
1066 5 enum DeviceType {
1067 6     UNDEFINED = 0;
1068 7     FAN = 1;
1069 8     PMP = 2;
1070 9     FCU = 3;
1071 10    VAV = 4;
1072 11    DH = 5;
1073 12    AHU = 6;
1074 13    BLR = 7;
1075 14    CDWS = 8;
1076 15    CH = 9;
1077 16    CHWS = 10;
1078 17    CT = 11;
1079 18    DC = 12;
1080 19    DFR = 13;
1081 20    DMP = 14;
1082 21    HWS = 15;
1083 22    HX = 16;

```

```

1080 22     MAU = 17;
1081 23     SDC = 18;
1082 24     UH = 19;
1083 25     PWR = 20;
1084 26     GAS = 21;
1085 27     AC = 22;
1086 28     OTHER = 23;
1087 29 }
1088 30
1088 31
1089 32 enum ValueType {
1090 33     VALUE_TYPE_UNDEFINED = 0;
1091 34     VALUE_CONTINUOUS = 1;
1092 35     VALUE_INTEGER = 2;
1093 36     VALUE_CATEGORICAL = 3;
1094 37     VALUE_BINARY = 4;
1095 38 }
1096 39
1096 40
1096 41 // Unique device identifier.
1097 42 string device_id = 1;
1098 43 // If applicable, the zone associated with the device (like VAVs).
1099 44 string namespace = 2;
1100 45 string code = 3;
1101 46 string zone_id = 4;
1102 47
1102 48
1103 49 // The type of device, VAV, AHU, etc.
1104 50 DeviceType device_type = 5;
1105 51 // Map of measurement name exposed by the device to the value type.
1106 52 map<string, ValueType> observable_fields = 6;
1107 53 // Map of setpoint name exposed by the device to their value type.
1108 54 map<string, ValueType> action_fields = 7;
1109 55 }

```

1109  
1110

#### B.1.4 DEVICEINFO EXAMPLE

1112

```

1113 1 device_id: "202194278473007104"
1114 2 namespace: "PHRED"
1115 3 code: "US-BLDG:AHU:AC-2"
1116 4 device_type: AHU
1117 5 observable_fields {
1118 6   key: "building_air_static_pressure_sensor"
1119 7   value: VALUE_CONTINUOUS
1120 8 }
1120 9 observable_fields {
1121 10  key: "building_air_static_pressure_setpoint"
1122 11  value: VALUE_CONTINUOUS
1123 12 }
1123 13 action_fields {
1124 14  key: "building_air_static_pressure_setpoint"
1125 15  value: VALUE_CONTINUOUS
1126 16 }
1127 17 action_fields {
1128 18  key: "cooling_percentage_command"
1129 19  value: VALUE_CONTINUOUS
1130 20 }
1130 21 action_fields {
1131 22  key: "exhaust_air_damper_percentage_command"
1132 23  value: VALUE_CONTINUOUS
1133 24 }

```

## 1134 B.2 OBSERVATION DATA PROTOS

1135  
1136 This includes the measurements from all devices in the building (VAV’s zone air temperature, gas  
1137 meter’s flow rate, etc.), provided at each time step. There are two proto definitions:

- 1138 1. **ObservationRequest**
- 1139 2. **ObservationResponse**

1140  
1141 To acquire the latest building state, at each timestep the building accepts  
1142 an `ObservationRequest` and returns an `ObservationResponse`. The  
1143 `ObservationRequest` contains a UTC timestamp of the requested observation, and list  
1144 of `SingleObservationRequests`. Each `SingleObservationRequest` is a tuple  
1145 of the `device_id` and the `measurement_name` that must match with a device and an  
1146 `observable_field` in one of the `DeviceInfos` exposed by the building. The building  
1147 returns an `ObservationResponse` that contains the UTC timestamp from the building, the  
1148 original `ObservationRequest`, and a list of `SingleObservationResponses`. Each  
1149 `SingleObservationResponse` contains the associated `SingleObservationRequest`,  
1150 the validity time of the measurement/observation, a boolean validity indicator, and the observation,  
1151 in native units, as a continuous, integer, categorical, binary or string value.

### 1152 B.2.1 OBSERVATIONREQUEST DEFINITION

```
1154 1 // Agent's request to get the current observation vector.
1155 2 message ObservationRequest {
1156 3   // UTC timestamp when the agent generated the request.
1157 4   google.protobuf.Timestamp timestamp = 1;
1158 5   // One or more individual requests.
1159 6   repeated SingleObservationRequest single_observation_requests = 2;
1160 7 }
1161 8
1162 9
1163 10 // A request to get a single measurement from a specific sensor.
1164 11 message SingleObservationRequest {
1165 12   // Unique device identifier.
1166 13   string device_id = 1;
1167 14   // Name of the sensor, e.g., zone_air_temperature.
1168 15   string measurement_name = 2;
1169 16 }
```

### 1169 B.2.2 OBSERVATIONREQUEST EXAMPLE

```
1170
1171 1 timestamp {
1172 2   seconds: 1682649309
1173 3   nanos: 942662000
1174 4 }
1175 5 single_observation_requests {
1176 6   device_id: "202194278473007104"
1177 7   measurement_name: "supply_fan_speed_frequency_sensor"
1178 8 }
1179 9 single_observation_requests {
1180 10  device_id: "202194278473007104"
1181 11  measurement_name: "mixed_air_temperature_sensor"
1182 12 }
1183 13 single_observation_requests {
1184 14  device_id: "202194278473007104"
1185 15  measurement_name: "outside_air_flowrate_setpoint"
1186 16 }
1187 17 single_observation_requests {
1188 18  device_id: "202194278473007104"
1189 19  measurement_name: "supply_air_temperature_sensor"
1190 20 }
```

### 1188 B.2.3 OBSERVATIONRESPONSE DEFINITION

```

1189
1190 1 // Building's response to an observation request message.
1191 2 message ObservationResponse {
1192 3   google.protobuf.Timestamp timestamp = 1;
1193 4   ObservationRequest request = 2;
1194 5   repeated SingleObservationResponse single_observation_responses = 3;
1195 6 }
1196 7
1197 8
1198 9
1199 10
1198 11 // Response for a single observation request.
1199 12 message SingleObservationResponse {
1200 13 // The validity time in UTC of the measurement.
1201 14 google.protobuf.Timestamp timestamp = 1;
1202 15 // Original request.
1203 16 SingleObservationRequest single_observation_request = 2;
1204 17 // Validity flag on the observation.
1205 18 bool observation_valid = 3;
1206 19 // Actual observed/measured value.
1207 20 oneof observation_value {
1208 21   float continuous_value = 4;
1209 22   int32 integer_value = 5;
1210 23   string categorical_value = 6;
1211 24   bool binary_value = 7;
1212 25   string string_value = 8;
1213 26 }
1214 27 }

```

### 1213 B.2.4 OBSERVATIONRESPONSE EXAMPLE

```

1214
1215
1216 1 timestamp {
1217 2   seconds: 1681110000
1218 3 }
1219 4 request {
1220 5   timestamp {
1221 6     seconds: 1682649309
1222 7     nanos: 942662000
1223 8   }
1224 9   single_observation_requests {
1225 10     device_id: "202194278473007104"
1226 11     measurement_name: "supply_fan_speed_frequency_sensor"
1227 12   }
1228 13   single_observation_requests {
1229 14     device_id: "202194278473007104"
1230 15     measurement_name: "mixed_air_temperature_sensor"
1231 16   }
1232 17   single_observation_requests {
1233 18     device_id: "202194278473007104"
1234 19     measurement_name: "outside_air_flowrate_setpoint"
1235 20   }
1236 21   single_observation_responses {
1237 22     timestamp {
1238 23       seconds: 1681109783
1239 24       nanos: 299000000
1240 25     }
1241 26     single_observation_request {
1242 27       device_id: "202194278473007104"
1243 28       measurement_name: "supply_fan_speed_frequency_sensor"
1244 29     }
1245 30     observation_valid: true
1246 31     continuous_value: 0.0
1247 32 }

```

```

1242 33 single_observation_responses {
1243 34   timestamp {
1244 35     seconds: 1681109783
1245 36     nanos: 299000000
1246 37   }
1247 38   single_observation_request {
1248 39     device_id: "202194278473007104"
1249 40     measurement_name: "mixed_air_temperature_sensor"
1250 41   }
1251 42   observation_valid: true
1252 43   continuous_value: 290.3909912109375
1253 44 }
1254 45 single_observation_responses {
1255 46   timestamp {
1256 47     seconds: 1681109783
1257 48     nanos: 299000000
1258 49   }
1259 50   single_observation_request {
1260 51     device_id: "202194278473007104"
1261 52     measurement_name: "outside_air_flowrate_setpoint"
1262 53   }
1263 54   observation_valid: true
1264 55   continuous_value: 8.825417518615723
1265 56 }

```

### 1264 B.3 ACTION DATA PROTOS

1266 This consists of the device setpoint values that the agent wants to set, provided at each timestep.  
1267 There are two relevant protos:

- 1269 1. **ActionRequest**
- 1270 2. **ActionResponse**

1272 The Environment converts the action from the agent into an ActionRequest and sends  
1273 it to the building. The building applies the request and returns an ActionResponse.  
1274 The ActionRequest contains the UTC timestamp from the Environment, and a list of  
1275 SingleActionRequests, one for each setpoint in the agent's action space. Each  
1276 SingleActionRequest contains a tuple of the device\_id, setpoint\_name, and re-  
1277 quested setpoint\_value, in native units. The device\_id must match with one of the  
1278 device\_ids in the DeviceInfos, and the setpoint\_name must match with one of the  
1279 action\_fields of the associated device. The ActionResponse contains the building's UTC  
1280 timestamp, the original ActionRequest, and a list of SingleActionResponses, one as-  
1281 sociated with each SingleActionRequest. The SingleActionResponse contains the  
1282 associated SingleActionRequest, a response type enumeration, and a string for additional  
1283 information.

#### 1284 B.3.1 ACTIONREQUEST DEFINITION

```

1285
1286 1 // Agent's request to the building with an action.
1287 2 message ActionRequest {
1288 3   // The UTC timestamp that the agent initiated the request.
1289 4   google.protobuf.Timestamp timestamp = 1;
1290 5   // One or more action requests to be performed.
1291 6   repeated SingleActionRequest single_action_requests = 2;
1292 7 }
1293 8
1294 9 // An action request to assign a value to one setpoint on one device.
1295 10 message SingleActionRequest {
1296 11   // The device being commanded.
1297 12   string device_id = 1;
1298 13   // Actual setpoint to be changed, like zone_air_temperature_setpoint.

```

```

1296 14 string setpoint_name = 2;
1297 15 oneof setpoint_value {
1298 16     float continuous_value = 3;
1299 17     int32 integer_value = 4;
1300 18     string categorical_value = 5;
1301 19     bool binary_value = 6;
1302 20     string string_value = 7;
1303 21 }
1304 22 }

```

1305

### 1306 B.3.2 ACTIONREQUEST EXAMPLE

1307

```

1308 1 timestamp {
1309 2     seconds: 1682649309
1310 3     nanos: 942662000
1311 4 }
1312 5 single_action_requests {
1313 6     device_id: "12945159110931775488"
1314 7     setpoint_name: "supply_air_static_pressure_setpoint"
1315 8     continuous_value: 186.8100128173828
1316 9 }
1317 10 single_action_requests {
1318 11     device_id: "12945159110931775488"
1319 12     setpoint_name: "supply_air_temperature_setpoint"
1320 13     continuous_value: 294.2592468261719
1321 14 }
1322 15 single_action_requests {
1323 16     device_id: "13761436543392677888"
1324 17     setpoint_name: "supply_water_temperature_setpoint"
1325 18     continuous_value: 310.9259338378906
1326 19 }
1327 20 single_action_requests {
1328 21     device_id: "13761436543392677888"
1329 22     setpoint_name: "differential_pressure_setpoint"
1330 23     continuous_value: 82737.09375
1331 24 }
1332 25 single_action_requests {
1333 26     device_id: "12945159110931775488"
1334 27     setpoint_name: "supervisor_run_command"
1335 28     continuous_value: -1.0
1336 29 }
1337 30 single_action_requests {
1338 31     device_id: "14409954889734029312"
1339 32     setpoint_name: "supervisor_run_command"
1340 33     continuous_value: -1.0
1341 34 }

```

1337

### 1338 B.3.3 ACTIONRESPONSE DEFINITION

1339

```

1340 1 // Building's response to an action request.
1341 2 message ActionResponse {
1342 3     // UTC timestamp of the building's response.
1343 4     google.protobuf.Timestamp timestamp = 1;
1344 5     // Original action request.
1345 6     ActionRequest request = 2;
1346 7     // Individual responses for each action.
1347 8     repeated SingleActionResponse single_action_responses = 3;
1348 9 }
1349 10
1350 11
1351 12
1352 13 // Building's response to a single action request.

```

```

1350 14 message SingleActionResponse {
1351 15   enum ActionResponseType {
1352 16     UNDEFINED = 0;
1353 17     // The building accepted the action as requested.
1354 18     ACCEPTED = 1;
1355 19     // The building is processing the request, but has not completed.
1356 20     PENDING = 2;
1357 21     // The action request timed out by request handler.
1358 22     TIMED_OUT = 3;
1359 23     // Request is rejected because the set value is not in an acceptable←
1360 24     range.
1361 25     REJECTED_INVALID_SETTING = 4;
1362 26     // Rejected because the setting is not enabled or available for ↔
1363 27     control.
1364 28     REJECTED_NOT_ENABLED_OR_AVAILABLE = 5;
1365 29     // A technician or control function overrode the action.
1366 30     REJECTED_OVERRIDE = 6;
1367 31     // The action was assigned to a device that does not exist.
1368 32     REJECTED_INVALID_DEVICE = 7;
1369 33     // The action was assigned to a valid device that's offline.
1370 34     REJECTED_DEVICE_OFFLINE = 8;
1371 35     UNKNOWN = 9;
1372 36     OTHER = 10;
1373 37   }
1374 38   SingleActionRequest request = 1;
1375 39   ActionResponseType response_type = 2;
1376 40   // Additional optional information related to the action/response.
1377 41   string additional_info = 3;
1378 42 }

```

### 1378 B.3.4 ACTIONRESPONSE EXAMPLE

```

1379 1 timestamp {
1380 2   seconds: 1681110000
1381 3 }
1382 4 request {
1383 5   timestamp {
1384 6     seconds: 1682649309
1385 7     nanos: 942662000
1386 8   }
1387 9   single_action_requests {
1388 10    device_id: "12945159110931775488"
1389 11    setpoint_name: "supply_air_static_pressure_setpoint"
1390 12    continuous_value: 186.8100128173828
1391 13  }
1392 14  single_action_requests {
1393 15    device_id: "12945159110931775488"
1394 16    setpoint_name: "supply_air_temperature_setpoint"
1395 17    continuous_value: 294.2592468261719
1396 18  }
1397 19  single_action_requests {
1398 20    device_id: "13761436543392677888"
1399 21    setpoint_name: "supply_water_temperature_setpoint"
1400 22    continuous_value: 310.9259338378906
1401 23  }
1402 24  single_action_responses {
1403 25    request {
1404 26      device_id: "12945159110931775488"
1405 27      setpoint_name: "supply_air_static_pressure_setpoint"
1406 28      continuous_value: 186.8100128173828
1407 29    }
1408 30    response_type: ACCEPTED

```

```

1404 31   additional_info: "2023-04-10 06:56:23.299000+00:00 129451591109317754↵
1405 32   88"
1406 33 }
1407 34 single_action_responses {
1408 35   request {
1409 36     device_id: "12945159110931775488"
1410 37     setpoint_name: "supply_air_temperature_setpoint"
1411 38     continuous_value: 294.2592468261719
1412 39   }
1413 40   response_type: ACCEPTED
1414 41   additional_info: "2023-04-10 06:56:23.299000+00:00 129451591109317754↵
1415 42   88"
1416 43 }
1417 44 single_action_responses {
1418 45   request {
1419 46     device_id: "13761436543392677888"
1420 47     setpoint_name: "supply_water_temperature_setpoint"
1421 48     continuous_value: 310.9259338378906
1422 49   }
1423 50   response_type: ACCEPTED
1424 51   additional_info: "2023-04-10 06:55:33.394000+00:00 137614365433926778↵
1425 52   88"

```

#### B.4 REWARD DATA PROTOS

This includes information used to calculate the reward, as expressed in cost in dollars, carbon footprint, and comfort level of occupants, provided at each time step. The Reward protos define the input and output messages for our 3C reward function (Cost Carbon and Comfort), which contains the code that converts them into a single scalar value, a requirement for most RL algorithms. There are two relevant protos:

1. **RewardInfo:** The values that are used as inputs to calculate the reward
2. **RewardResponse:** Containing the scalar reward signal obtained by passing the above functions into our 3C reward function

The building updates the `RewardInfo` at each timestep and provides the reward function necessary inputs to compute the 3C Reward Function. The data contained in the `RewardInfo` is bounded by the step's interval from `start_timestamp` to `end_timestamp` in UTC. The `RewardInfo` has mean energy rate estimates (i.e. power in Watts) that can be treated as constants over the interval. Given the interval and a constant rate value over the interval, the reported power in Watts can be easily converted into energy in kWh. The `RewardInfo` contains maps of three types of specialized data structures:

- The `ZoneRewardInfo` message provides information about the zone air temperature measurements, temperature setpoints, airflow rate and setpoint, and average occupancy for the time step. Each instance is indexed by its unique zone ID.
- The `AirHandlerRewardInfo` message describes the combined electrical power in W use of the intake/exhaust blowers, and the electrical power in W of the refrigeration cycle. Since a building may have more than one air handler, the air handler objects are values in a map keyed by the air handlers' device IDs.
- The `BoilerRewardInfo` contains the average electrical power in W used by the pumps to circulate water through the building, and the average natural gas power in W used to heat the water in the boiler. Since there may be more than one hot water cycle in the building, each `ZoneRewardInfo` is placed into a map keyed by the hot water device's ID.

The reward function converts the current `RewardInfo` into the `RewardResponse` for the same interval as the `RewardInfo`. The agent's reward signal is `agent_reward.value`. Since the reward returned to the agent is a function of multiple factors, it is useful for analysis to show the individual components, such as carbon mass emitted, and the electrical and gas costs for the step.

```

1458 B.4.1 REWARDINFO DEFINITION
1459
1460 1 message RewardInfo {
1461 2 // Information about each zone in the time step for computing reward.
1462 3 message ZoneRewardInfo {
1463 4 // Heating setpoint of the zone at the timestep in K.
1464 5 float heating_setpoint_temperature = 1;
1465 6
1466 7
1466 8 // Cooling setpoint of the zone at the timestep in K.
1467 9 float cooling_setpoint_temperature = 2;
1468 10
1469 11
1469 12 // Average zone air temperature measured in the zone in K.
1470 13 float zone_air_temperature = 3;
1471 14
1472 15
1472 16 // Setpoint for air flow ventilation in the zone in m^3/s.
1473 17 float air_flow_rate_setpoint = 4;
1474 18
1475 19
1476 20 // Actual ventilation air flow in the zone in m^3/s.
1477 21 float air_flow_rate = 5;
1478 22
1479 23
1479 24 // Average occupancy in the zone over the time step in number of
1480 25 // people in the zone.
1481 26 float average_occupancy = 6;
1482 27 }
1483 28
1483 29 // Information about the air handler energy consumption for computing ←
1484 30 reward.
1485 31 message AirHandlerRewardInfo {
1486 32 // Cumulative electrical power in W applied to blowers.
1487 33 float blower_electrical_energy_rate = 1;
1488 34
1488 35 // Cumulative electrical energy rate applied in W for air ←
1489 36 conditioning. This
1490 37 // represents the total power applied for running a refrigeration or
1491 38 // heat pump cycles (includes running a compressor and pumps to
1492 39 // recirculate refrigerant.).
1493 40 float air_conditioning_electrical_energy_rate = 2;
1494 41 }
1495 42
1496 43 // Information about the boiler that provides heated water for VAVs.
1497 44 message BoilerRewardInfo {
1498 45 // Energy rate consumed in W by natural gas for heating water.
1499 46 float natural_gas_heating_energy_rate = 1;
1500 47 // Cumulative electrical power in W for water recirculation pumps.
1501 48 float pump_electrical_energy_rate = 2;
1502 49 }
1503 50
1503 51 // Start and end timestamps bound the timestep of the reward ←
1504 52 information.
1505 53 google.protobuf.Timestamp start_timestamp = 1;
1506 54 google.protobuf.Timestamp end_timestamp = 2;
1507 55
1508 56 // Unique ID of the agent (controller). This should reflect the
1509 57 // attributes of the RL models, including the type of algo and its
1510 58 // parameters.
1511 59 string agent_id = 3;
1512 60

```

```

1512 61
1513 62 // Unique ID of the scenario being executed. This should reflect the ↵
1514 63 details
1515 64 // of the scenario. In simulation, it should identify the canonical ↵
1516 65 scenario.
1517 66 // In real world, it should define the building and start date/time.
1518 67 string scenario_id = 4;
1519 68
1520 69 // Map with zone_id and zone reward info for all zones in the building
1521 70 // under control of the agent. The zone_id could be a unique room ↵
1522 71 number,
1523 72 // or the specific zone coordinates: (i.e., 'z_i,z_j') from the ↵
1524 73 simulation.
1525 74 map<string, ZoneRewardInfo> zone_reward_infos = 5;
1526 75
1527 76 // Information about the air handlers' energy consumption required to
1528 77 // calculate the reward.
1529 78 map<string, AirHandlerRewardInfo> air_handler_reward_infos = 6;
1530 79
1531 80 // Information about the boilers' energy consumption required to ↵
1532 81 compute the
1533 82 // reward.
1534 83 map<string, BoilerRewardInfo> boiler_reward_infos = 7;
1535 84 }

```

1535

1536

#### 1537 B.4.2 REWARDINFO EXAMPLE

```

1538 1 start_timestamp {
1539 2   seconds: 1681109700
1540 3 }
1541 4 end_timestamp {
1542 5   seconds: 1681110000
1543 6 }
1544 7 agent_id: "baseline_policy"
1545 8 scenario_id: "baseline_collect"
1546 9 zone_reward_infos {
1547 10  key: "rooms/1000004614278"
1548 11  value {
1549 12    heating_setpoint_temperature: 289.0
1550 13    cooling_setpoint_temperature: 298.0
1551 14    zone_air_temperature: 293.5944519042969
1552 15    air_flow_rate_setpoint: 258.0
1553 16    air_flow_rate: 12.0
1554 17  }
1555 18 }
1556 19 zone_reward_infos {
1557 20  key: "rooms/1000004658174"
1558 21  value {
1559 22    heating_setpoint_temperature: 289.0
1560 23    cooling_setpoint_temperature: 298.0
1561 24    zone_air_temperature: 293.4277648925781
1562 25    air_flow_rate_setpoint: 60.0
1563 26  }
1564 27 }
1565 28 zone_reward_infos {
1566 29  key: "rooms/1000004658175"
1567 30  value {
1568 31    heating_setpoint_temperature: 289.0
1569 32    cooling_setpoint_temperature: 298.0
1570 33    zone_air_temperature: 293.03887939453125
1571 34    air_flow_rate_setpoint: 185.0

```

```

1566 35     air_flow_rate: 4.001242637634277
1567 36   }
1568 37 }
1569 38 zone_reward_infos {
1570 39   key: "rooms/1000004658176"
1571 40   value {
1572 41     heating_setpoint_temperature: 289.0
1573 42     cooling_setpoint_temperature: 298.0
1574 43     zone_air_temperature: 293.53887939453125
1575 44     air_flow_rate_setpoint: 145.0
1576 45     air_flow_rate: 53.0
1577 46   }
1578 47 }
1579 48 air_handler_reward_infos {
1580 49   key: "12945159110931775488"
1581 50   value {
1582 51   }
1583 52 }
1584 53 air_handler_reward_infos {
1585 54   key: "14409954889734029312"
1586 55   value {
1587 56   }
1588 57 }
1589 58 boiler_reward_infos {
1590 59   key: "13761436543392677888"
1591 60   value {
1592 61     pump_electrical_energy_rate: 1527.1470947265625
1593 62   }
1594 63 }

```

1591

### 1592 B.4.3 REWARDRESPONSE DEFINITION

1593

```

1594 1 // The return reward signal from the reward function. While the ←
1595 2     principal
1596 3 // signal is the agent reward and should be returned to the RL agent, ←
1597 4     the
1598 5 // other fields provide useful information for tracking and monitoring.
1599 6 // One EnergyRewardResponse is associated with each EnergyRewardInfo.
1600 7 message RewardResponse {
1601 8
1602 9     // Complete reward signal to be returned to the agent.
1603 10     float agent_reward_value = 1;
1604 11
1605 12     // Cumulative productivity is measured in USD, and represents the ←
1606 13     total
1607 14     // estimated productivity of the building.
1608 15     float productivity_reward = 2;
1609 16
1610 17     // Total electrical energy cost estimate in USD.
1611 18     float electricity_energy_cost = 3;
1612 19
1613 20
1614 21     // Total natural gas energy cost in USD.
1615 22     float natural_gas_energy_cost = 4;
1616 23
1617 24
1618 25     // Estimated carbon emitted in kg.
1619 26     float carbon_emitted = 5;
1620 27
1621 28
1622 29     // Estimated carbon cost in USD.

```

```

1620 30 float carbon_cost = 6;
1621 31
1622 32
1623 33 // Productivity weight parameter.
1624 34 float productivity_weight = 7;
1625 35
1626 36
1627 37 // Energy Cost Weight parameter.
1628 38 float energy_cost_weight = 8;
1629 39
1630 40
1631 41 // Carbon emission weight parameter.
1632 42 float carbon_emission_weight = 9;
1633 43
1634 44
1635 45 // Productivity factor (avg labor value of one person-hour).
1636 46 float person_productivity = 10;
1637 47
1638 48
1639 49 // Total average occupancy across all zones.
1640 50 float total_occupancy = 11;
1641 51
1642 52
1643 53 // Reward scale for normalizing the reward
1644 54 float reward_scale = 12;
1645 55
1646 56
1647 57 // Reward shift for normalizing the reward
1648 58 float reward_shift = 13;
1649 59
1650 60
1651 61 // Total productivity regret = max productivity - actual productivity
1652 62 float productivity_regret = 14;
1653 63
1654 64
1655 65 // Normalized productivity regret
1656 66 float normalized_productivity_regret = 15;
1657 67
1658 68
1659 69 // Normalized energy cost =
1660 70 // combined_energy_cost /
1661 71 // (max_electricity_energy_cost + max_natural_gas_energy_cost)
1662 72 float normalized_energy_cost = 16;
1663 73
1664 74
1665 75 // Normalized carbon emission =
1666 76 // combined_carbon_emission /
1667 77 // (max_electricity_carbon_emission + max_natural_gas_carbon_emission←
1668 78 )
1669 79 float normalized_carbon_emission = 17;
1670 80
1671 81 // Start and end timestamps bound the timestep of the reward ←
1672 82 // information.
1673 83 google.protobuf.Timestamp start_timestamp = 18;
1674 84 google.protobuf.Timestamp end_timestamp = 19;

```

#### 1669 B.4.4 REWARDRESPONSE EXAMPLE

```

1671 1 agent_reward_value: -0.00222194055095315
1672 2 electricity_energy_cost: 0.022907206788659096
1673 3 carbon_emitted: 0.011416268534958363
1674 4 productivity_weight: 0.5

```

```

1674 5 energy_cost_weight: 0.20000000298023224
1675 6 carbon_emission_weight: 0.30000001192092896
1676 7 person_productivity: 300.0
1677 8 reward_scale: 1.0
1678 9 normalized_energy_cost: 0.0090464623644948
1679 10 normalized_carbon_emission: 0.0013754934770986438
1680 11 start_timestamp {
1681 12   seconds: 1681109700
1681 13 }
1682 14 end_timestamp {
1683 15   seconds: 16811100

```

## C SIMULATOR DESIGN CONSIDERATION DETAILS

A simulator models the physical system dynamics of the building, devices, and external weather conditions, and can train the control agent interactively, if the following desiderata are achieved:

1. The simulation must produce the same observation dimensionality as the actual real building. In other words, each device-measurement present in the real building must also be present in the simulation.
2. The simulation must accept the same actions (device-setpoints) as the real building.
3. The simulation must return the reward input data described above (zone air temperatures, energy use, and carbon emission).
4. The simulation must propagate, estimate, and compute the thermal dynamics of the actual real building and generate a state update at each timestep.
5. The simulation must model the dynamics of the HVAC system in the building, including thermostat response, setpoints, boiler, air conditioning, water circulation, and air circulation. This includes altering the HVAC model in response to a setpoint change in an action request.
6. The time required to recalculate a timestep must be short enough to train a viable agent in a reasonable amount of time. For example, if a new agent should be trained in under three days (259,200 seconds), requiring 500,000 steps, the average time required to update the building should be 0.5 seconds or less.
7. The simulator must be configurable to a target building with minimal manual effort.

We believe our simulation system meets all of these listed requirements.

## D DERIVATION FOR TENSORIZED FINITE DIFFERENCE (FD) EQUATIONS

This appendix describes the method of calculating the flow of heat and the resulting temperatures throughout the building.

### D.1 ASSEMBLING THE ENERGY BALANCE

The fundamental energy balance for a general-purpose closed body is formulated in Equation 3. The first term represents the effects of non-stationary heat dissipation or heat absorption over time over volume of the body.  $Q$  represents the energy absorbed or released per unit volume and is a function of the mass and heat capacity of the body. The second term represents thermal flux over the surface of the body, where  $\mathbf{n}$  is the unit normal vector of the surface  $S$  and  $\mathbf{F}$  is the specific energy absorbed or released through the surface. Common modes of thermal flux include conduction, convection, and radiation. The right side of the equation represents the total energy absorbed by the body across the system boundary, or via an external source or sink.

$$\frac{d}{dt} \int_{V(t)} Q dV + \oint_{S(t)} \mathbf{n} \cdot \mathbf{F} dS = \int_{V(t)} P dV \quad (2)$$

1728 To enable computation, we divide the body into small discrete units, called **Control Volumes (CV)**,  
 1729 and iteratively calculate temperature on each on each CV using the method of Finite Differences  
 1730 (FD).

1731 We model three modes of heat transfer into each CV: forced convection, conduction, and external  
 1732 source.

1733 Forced convection  $Q^{conv}$  is based on energy exchange by moving air (or any other fluid, in general),  
 1734 and conduction,  $Q^{cond}$  is the exchange of energy through solid objects, such as walls. External  
 1735 sources (or sinks)  $Q^x$  represent the heating or cooling from external devices, such as electric heating  
 1736 coils, diffusers, etc.

1738 Each CV has the capacity to absorb heat over time, which is expressed as  $\frac{dU}{dt}$ , governed by its heat  
 1739 capacity,  $c$ .

1740 These factors allow us to construct an energy balance equation that conserves energy  $Q^{in} - Q^{out} =$   
 1741  $\frac{dU}{dt}$ .

1743 We assume that the ceilings and floors are adiabatic, fully insulated, not allowing any heat exchange.  
 1744 This reduces the problem to a 2D problem, with 3D control volumes that can only exchange energy  
 1745 laterally.

1746 Our FD objective is to solve for the temperature at each CV within the building, which presents  $N$   
 1747 unknowns and  $N$  equations, where  $N$  is the number of CVs in the FD grid.

1748 Rather than creating separate spacial cases in the FD equations for exterior, boundary, and interior  
 1749 CVs, we would like to create a single equation that can be computed across the entire grid. This  
 1750 equation can then be tensorized using the Tensorflow matrix library, and accelerated with GPUs or  
 1751 TPUs.

1752 We label each four interacting surfaces of the CV: left = 1, right = 3, bottom = 2, and top = 4.

1753 Then, for a discrete unit of time  $\Delta t$  we specify energy exchange across the surfaces as  
 1754  $Q_1, Q_2, Q_3, Q_4$  and adopt the arbitrary, but consistent convention that energy flows into surfaces  
 1755 1 and 2, and out of surfaces 3, and 4. (Of course, energy can flow the other direction too, but that  
 1756 will be indicates with a negative value.) Our convention also assumes that external energy flows into  
 1757 the CV.

1759 That allows us to construct the energy balance as:

$$1761 \quad Q^x + Q_1^{cond} + Q_1^{conv} + Q_2^{cond} + Q_2^{conv} - Q_3^{cond} - Q_3^{conv} - Q_4^{cond} - Q_4^{conv} = \frac{dU}{dt} \quad (3)$$

## 1766 D.2 COMPUTING HEAT TRANSFER VIA CONDUCTION, CONVECTION, AND THERMAL 1767 ABSORPTION

1768 We apply the Fourier’s Law of conduction, illustrated in Figure 8, which is the rate of transfer in  
 1769 Watts:

$$1773 \quad \dot{Q}^{cond} = -\frac{kA}{L} \frac{dT}{dt} \quad (4)$$

1777 Which is approximated over the discrete CV as:

$$1781 \quad \dot{Q}^{cond} \approx -\frac{kA}{L} \frac{\Delta T}{\Delta t} \quad (5)$$

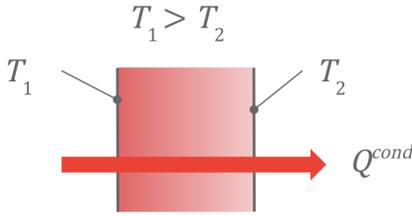


Figure 8: Conduction Heat Transfer

Where  $k$  is the thermal conductivity of the material,  $A$  is the flux area perpendicular to the flow of heat,  $L$  is the distance traveled through the material,  $\Delta T$  is the temperature difference in the source and sink, and  $\Delta t$  is a discrete time step interval.

We can remove the dot (time derivative) by multiplying by discrete unit time, and converting thermal power (energy per unit time) into energy:

$$Q^{cond} \approx -\frac{kA}{L} \frac{\Delta T}{\Delta t} \times 1 = -\frac{kA}{L} \Delta T \quad (6)$$

Let's orient the conductivity equation along the horizontal ( $u$ ) and the vertical directions ( $v$ ).

For the horizontal heat transfer:

$$Q_{1,3}^{cond} = -\frac{kvz}{u} \Delta T \quad (D.5) \quad (7)$$

And for vertical heat transfer:

$$Q_{2,4}^{cond} = -\frac{kuz}{v} \Delta T \quad (8)$$

Where  $z$  is the 3rd dimension size, which is the distance from the floor to the ceiling, and  $A = vz$  and  $A = uz$  for horizontal and vertical flux surface areas.

This is good for modeling heat exchange through solid objects, but we also need to model the heat exchanges from the outside across the boundary to the interior via forced air convection (i.e., wind).

For convection, we'll apply Newton's Law of Cooling, illustrated in Figure 9 for modeling heat transfer via forced air currents across a surface  $A$ , perpendicular to the flow of heat as:

$$Q^{cond} = -hA\Delta T \quad (9)$$

The negative sign in Equations 4 - 9 are due to the fact that energy flows in the direction opposite of the temperature gradient,  $\Delta T$ , i.e., from high to low.

Here,  $h$  is the convection coefficient and is a function of the amount of air blowing over the exterior surface of the wall.

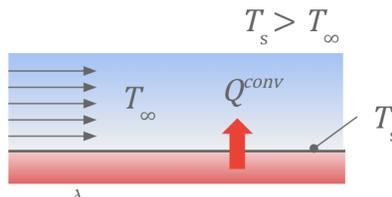


Figure 9: Convection Heat Transfer

We define the three types of CVs:

1. **Exterior CVs** are CVs that represent the ambient weather conditions, such as  $T_\infty$ , which are not calculated by the FD calculator, just specified by the current input conditions.
2. **Interior CVs** are CVs where all four sides are adjacent to non-exterior CVs (Figure 10).
3. **Boundary CVs** are CVs that share one or two faces with exterior CVs and one two or three faces with interior CVs. These CVs require special handling, since they represent the transfer of energy between the outside and the inside of the building. Boundary CVs that share two sides with the exterior are **Corner CVs** (Figure 11) and boundary CVs that share only one side with an exterior CV are **Edge CVs** (Figure 12).

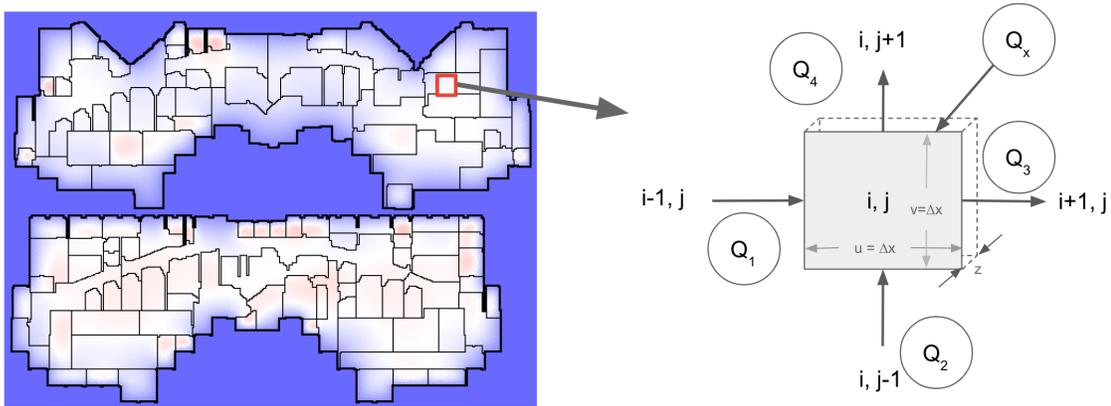


Figure 10: Interior Control Volumes

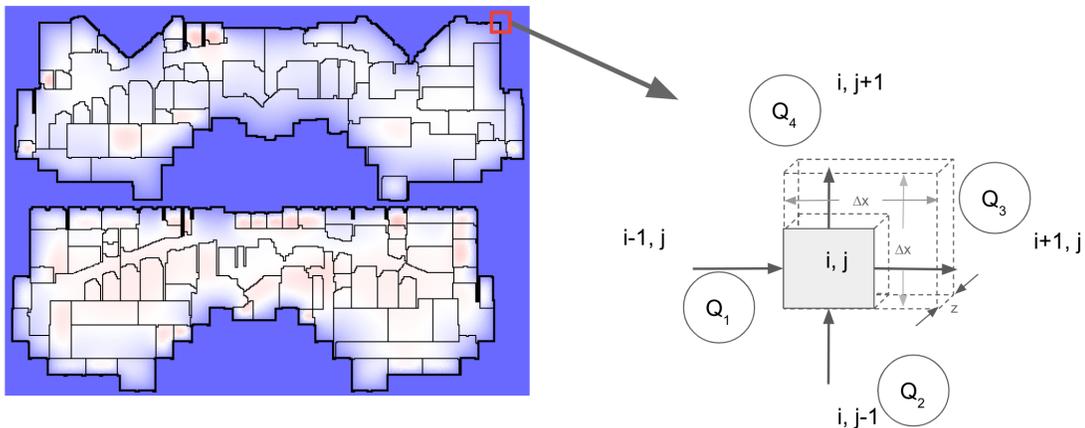


Figure 11: Boundary Corner Control Volumes

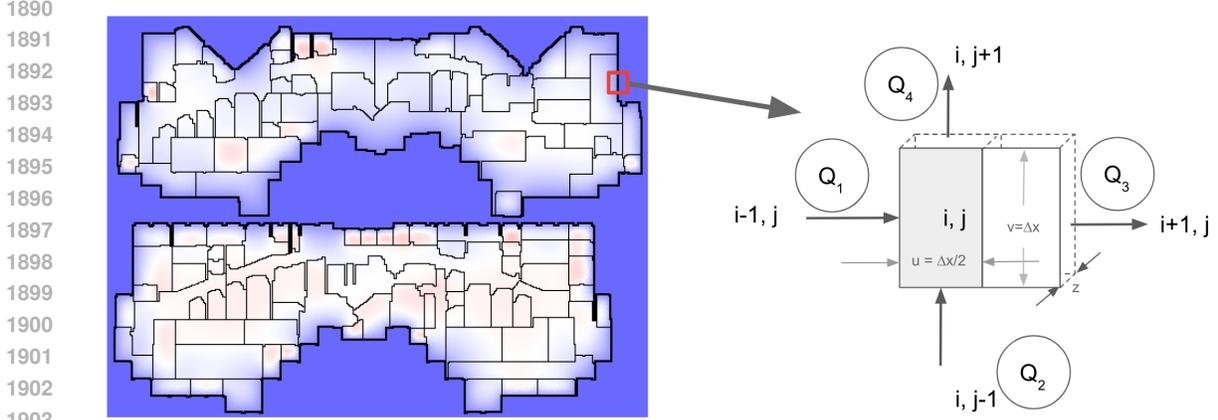


Figure 12: Boundary Edge Control Volumes

1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913

The temperatures that are estimated in FD represent the center of the control volume, or its mean. In the case of convection, the temperatures at the exterior surface of the wall are unknown and have to be calculated. Therefore, the center of the Edge CV represents the surface temperature and is split halfway between the outside and inside, where the volume of an edge CV is half of the mass of an interior CV. Similarly, a corner CV is cut in half in both directions, and is one quarter the volume of an interior CV.

1914  
1915

Since we are assuming rectangular CVs, note that  $v = v_1 = v_3$ , and  $u = u_2 = u_4$ .

1916  
1917  
1918  
1919

Since outside temperatures and HVAC responses vary, we have a non-stationary thermal system where the flow of energy through the CVs that is not constant. This requires us to evaluate the right-hand term in Equation 3 that allows the volume to absorb or dissipate heat over time, which is governed by the mass  $m = \rho V = \rho uvz$ , heat capacity  $c$  and rate of change of temperature  $\frac{dT}{dt}$ .

1920  
1921  
1922  
1923

$$\frac{dU}{dt} = cm \frac{dT}{dt} = c\rho V \frac{dT}{dt} = c\rho uvz \frac{dT}{dt} \quad (10)$$

1924  
1925

Equation 10 can be approximated over the small differential CV as:

1926  
1927  
1928  
1929

$$\frac{dU}{dt} \approx c\rho uvz \frac{T_{i,j} - T_{i,j}^{(-)}}{\Delta t} \quad (11)$$

1930  
1931  
1932  
1933

where  $T_{i,j}^{(-)}$  is the temperature if the  $i, j$  CV at the previous time step and the time step interval is  $\Delta t$ , which can be treated as a fixed parameter.

### 1934 1935 D.3 SOLVING FOR THE TEMPERATURE AT EACH CV

1936  
1937  
1938  
1939  
1940  
1941

To enable accelerating the calculation using tensor operations, we would like to define a single equation for all CV that do not require (a) conditionals, (b) for loops, or (c) referencing neighboring CVs. That objective will require the construction of a few auxiliary matrices, and every CV will have convection and conduction components that may be disabled with zero-valued convection and conduction coefficients as appropriate.

1942  
1943

Combining the Energy Balance in Equation 4 with the conduction and convection equations (Equations 7-10) we can include all terms for all faces on the  $i, j$  CV. Our goal is to solve for  $T_{i,j}$  which can then be run over multiple sweeps to convergence.

$$\begin{aligned}
Q_x - k_1 v z \frac{T_{i,j} - T_{i-1,j}}{u} - h_1 v z (T_{i,j} - T_\infty) - k_2 u z \frac{T_{i,j} - T_{i,j-1}}{v_2} - h_2 v z (T_{i,j} - T_\infty) + \\
+k_3 v z \frac{T_{i+1,j} - T_{i,j}}{u_3} + h_3 v z (T_\infty - T_{i,j}) + k_4 u z \frac{T_{i,j+1} - T_{i,j}}{v_4} + h_4 v z (T_\infty - T_{i,j}) = \quad (12) \\
= \frac{c \rho u v z}{\Delta t} (T_{i,j} - T_{i,j}^{(-)})
\end{aligned}$$

Next, we want to solve for temperature  $T_{i,j}$  by rearranging the terms, which provides a single equation that can be used to calculate CV temperatures for both boundary and interior CVs.

$$T_{i,j} = \frac{Q_x + v z \left[ \frac{k_1}{u} T_{i-1,j} + h_1 T_\infty + \frac{k_3}{u} T_{i+1,j} + h_3 T_\infty \right] + u z \left[ \frac{k_2}{v} T_{i,j-1} + h_2 T_\infty + \frac{k_4}{v} T_{i,j+1} + h_4 T_\infty \right] + \frac{c \rho u v z}{\Delta t} T_{i,j}^{(-)}}{v z \left[ \frac{k_1}{u} + h_1 + \frac{k_3}{u} + h_3 \right] + u z \left[ \frac{k_2}{v} + h_2 + \frac{k_4}{v} + h_4 \right] + \frac{c \rho u v z}{\Delta t}} \quad (13)$$

#### D.4 TENSORIZING THE TEMPERATURE ESTIMATE

Equation 13 can be used iterative, but to exploit the acceleration from matrix operations on GPUs and TPUs using the TensorFlow Library, we'll want to reshape the equation slightly for a single tensor pipeline that doesn't iterate over individual CVs.

Furthermore, we can avoid referencing neighboring temperatures ( $T_{i-1,j}, T_{i+1,j}, T_{i,j-1}, T_{i,j+1}$ ) in the pipeline by creating four \*shifted\* temperature Tensors,  $T_1 = \text{shift}(T, 3)$ ,  $T_3 = \text{shift}(T, \text{LEFT})$ ,  $T_2 = \text{shift}(T, \text{UP})$ ,  $T_4 = \text{shift}(T, \text{DOWN})$ .

We can also frame oriented conductivity as a Tensors left  $K_1$ , right  $K_3$ , below  $K_2$ , above  $K_4$ , where:

$$k_{1,i,j} = \begin{cases} k_{i,j} & \text{CVs at } i, j \text{ and } i-1, j \text{ are interior or boundary} \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

$$k_{3,i,j} = \begin{cases} k_{i,j} & \text{CVs at } i, j \text{ and } i+1, j \text{ are interior or boundary} \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

$$k_{2,i,j} = \begin{cases} k_{i,j} & \text{CVs at } i, j \text{ and } i, j-1 \text{ are interior or boundary} \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

$$k_{4,i,j} = \begin{cases} k_{i,j} & \text{CVs at } i, j \text{ and } i, j+1 \text{ are interior or boundary} \\ 0 & \text{otherwise} \end{cases} \quad (17)$$

Note that the conductivity matrix  $K$  is a fixed input parameter for the building.

Applying the same reasoning, we can generate four oriented convection Tensors,  $H_1, H_2, H_3, H_4$  as:

$$h_{1,i,j} = \begin{cases} h & \text{CV at } i, j \text{ is boundary and CV at } i-1, j \text{ is exterior} \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

$$h_{3,i,j} = \begin{cases} h & \text{CV at } i, j \text{ is boundary and CV at } i+1, j \text{ is exterior} \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

$$h_{2,i,j} = \begin{cases} h & \text{CV at } i, j \text{ is boundary and CV at } i, j+1 \text{ is exterior} \\ 0 & \text{otherwise} \end{cases} \quad (20)$$

1998  
 1999  
 2000  
 2001  
 2002  
 2003  
 2004  
 2005  
 2006  
 2007  
 2008  
 2009  
 2010  
 2011  
 2012  
 2013  
 2014  
 2015  
 2016  
 2017  
 2018  
 2019  
 2020  
 2021  
 2022  
 2023  
 2024  
 2025  
 2026  
 2027  
 2028  
 2029  
 2030  
 2031  
 2032  
 2033  
 2034  
 2035  
 2036  
 2037  
 2038  
 2039  
 2040  
 2041  
 2042  
 2043  
 2044  
 2045  
 2046  
 2047  
 2048  
 2049  
 2050  
 2051

$$h_{4,i,j} = \begin{cases} h & \text{CV at } i,j \text{ is boundary and CV at } i,j-1 \text{ is exterior} \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

Note that  $h$  is a time-dependent constant that represents the amount of airflow over the surface of the building, assumed to be uniformly applied on all exterior walls of the building.

Finally, we classify each boundary CV as TOP-LEFT CORNER, TOP-RIGHT CORNER, BOTTOM-LEFT CORNER, BOTTOM-RIGHT CORNER or LEFT EDGE, RIGHT EDGE, TOP EDGE, or BOTTOM EDGE in order to form Tensors  $U$  and  $V$ , which are the CV widths and heights.

$$u_{i,j} = \begin{cases} \frac{\Delta x}{2} & \text{CV at } i,j \text{ is BOUNDARY and ANY CORNER or TOP or BOTTOM EDGE} \\ \Delta x & \text{otherwise} \end{cases} \quad (22)$$

$$v_{i,j} = \begin{cases} \frac{\Delta x}{2} & \text{CV at } i,j \text{ is BOUNDARY and ANY CORNER or LEFT or RIGHT EDGE} \\ \Delta x & \text{otherwise} \end{cases} \quad (23)$$

where  $\Delta x$  is the fixed horizontal and vertical dimension of an INTERIOR CV.

Now we can complete the Tensor expression of the FD equation:

$$T = \left[ Q_x + Vz [K_1 U^{-1} T_1 + H_1 T_\infty + K_3 U^{-1} T_3 + H_3 T_\infty] \right. \\ \left. + Uz [K_2 V^{-1} T_2 + H_2 T_\infty + K_4 V^{-1} T_4 + H_4 T_\infty] \right. \\ \left. + \frac{CPUVz}{\Delta t} T^{(-)} \right] \\ \cdot \left[ Vz [K_1 U^{-1} + H_1 + K_3 U^{-1} + H_3] + Uz [K_2 V^{-1} + H_2 + K_4 V^{-1} + H_4] + \frac{CPUVz}{\Delta t} \right]^{-1} \quad (24)$$

For each timestep, we execute Equation 24 as single-step tensor operations until convergence, where the maximum change across all CVs between current and last iteration is less than a conservative lower threshold,  $\epsilon \leq 0.01^\circ C$

## E SIMULATOR CONFIGURATION PROCEDURE DETAILS

To configure the simulator, we require two type of information on the building:

1. Floorplan blueprints. This includes the size and shapes of rooms and walls for each floor.
2. HVAC metadata. This includes each device, its name, location, setpoints, fixed parameters and purpose.

We preprocess the detailed floorplan blueprints of the building, and extract a grid that gives us an approximate placement of walls and how rooms are divided. This is done via the following procedure:

1. Using threshold  $t$ , binarize the floorplan image into a grid of 0s and 1s.
2. Find and replace any large features that need to be removed (such as doors, a compass, etc)
3. Iteratively apply standard binary morphology operations (erosion and dilation) to the image to remove noise from background, while preserving the walls.

4. Resize the image, such that each pixel represents exactly one control volume
5. Run a connected components search to determine which control volumes are exterior to the building, and mark them accordingly
6. Run a DFS over the grid, and reduce every wall we encounter to be only a single control volume thick in the case of interior wall, and double for exterior wall

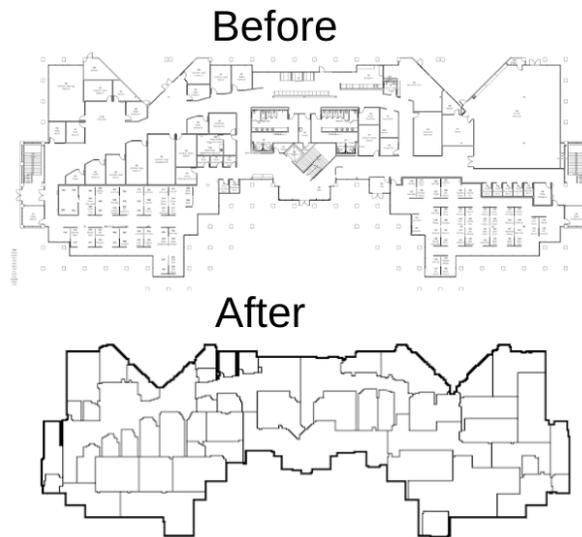


Figure 13: Before and after images of the floorplan preprocessing algorithm

We also employ a simple user interface to label the location of each HVAC device on the floorplan grid. This information is passed into our simulator, and a custom simulator for the new building, with roughly accurate HVAC and floor layout information, is created. This allows us to then calibrate this simulator using the real world data, which will now match the simulator in terms of device names and locations.

We tested this pipeline on SB1, which consisted of two floors with combined surface area of 93,858 square feet, and has 127 HVAC devices. Given floorplans and HVAC layout information, a single technician was able to generate a fully specified simulation in under three hours. This customized simulator matched the real building in every device, room, and structure.

## F CALIBRATION HYPERPARAMETER TUNING DETAILS

The hyperparameter tuning was performed over a seven day period on 200 CPUs.

Table 4: Thermal properties that were set by the calibration process, with min/max bounds and selected values.

HYPERPARAMETER	MIN	MAX	BEST
CONVECTION_COEFFICIENT ( $W/m^2/K$ )	5	800	357
EXTERIOR_CV_CONDUCTIVITY ( $W/m/K$ )	0.01	1	0.83
EXTERIOR_CV_DENSITY ( $kg/m^3$ )	0	3000	2359
EXTERIOR_CV_HEAT_CAPACITY ( $J/Kg/K$ )	100	2500	2499
INTERIOR_WALL_CV_CONDUCTIVITY ( $W/m/K$ )	5	800	5
INTERIOR_WALL_CV_DENSITY ( $kg/m^3$ )	0.5	1500	1500
INTERIOR_WALL_CV_HEAT_CAPACITY ( $J/Kg/K$ )	500	1500	1499
SWAP_PROB	0	1	0.003
SWAP_RADIUS	0	50	50

## G ADDITIONAL SPATIAL ERROR VISUALIZATIONS

Here we present some other visuals that may be enlightening.

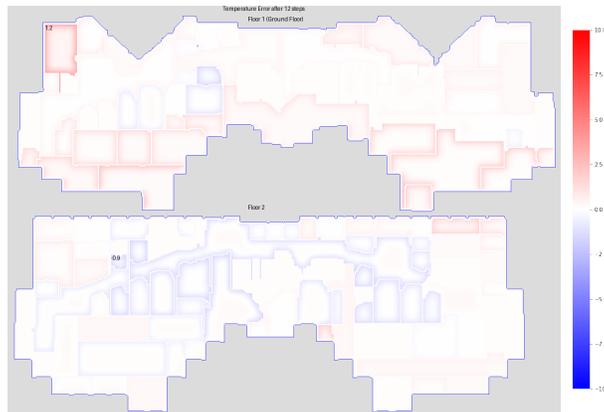


Figure 14: Visualization of simulator drift after only a single hour, on the validation data. As can be clearly seen, at this point there is almost no error.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173

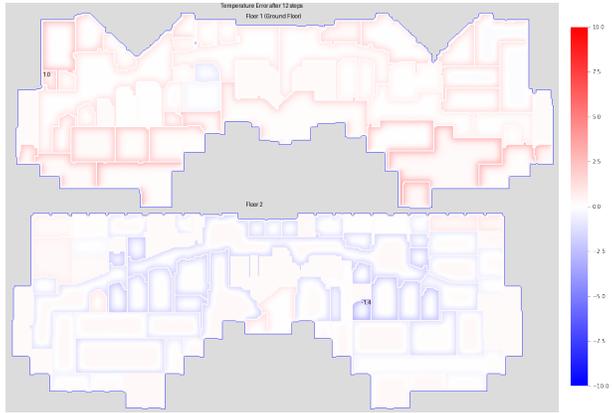


Figure 15: Visualization of simulator drift after only a single hour, on the train data. Again, there is almost no error.

2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192

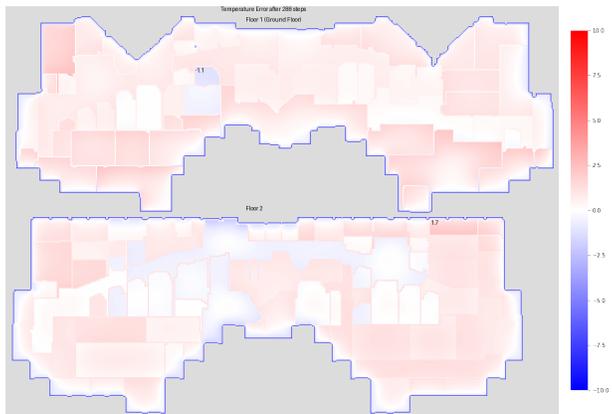


Figure 16: Visualization of simulator drift after one day, on the train data.

2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211

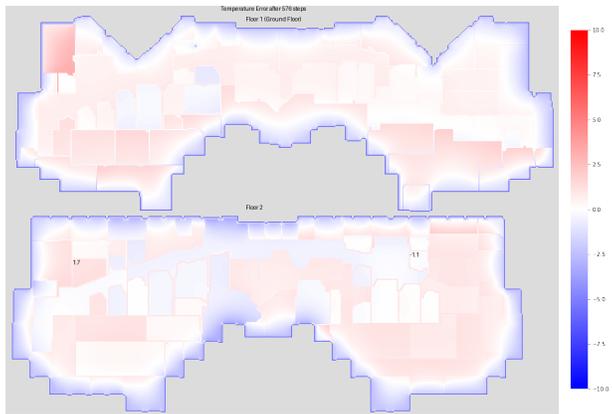


Figure 17: Visualization of simulator drift after two days, on the train data. Interestingly, this looks better than it did after only one day.

## H SIMULATOR SAC AGENT TRAINING DETAILS AND PERFORMANCE ANALYSIS

We will now go into more details on the simulator SAC agent training and performance as compared to the baseline.

Each agent was trained on a single CPU, with the entire training session lasting 6 days. We restricted the action space to supply air and water temperature setpoints. For the observation space, we found that providing the agent with the dozens of temperature sensors was too much noisy information and not useful. Instead, we provided the agent with a histogram, grouping temperatures into  $1^\circ$  Celsius bins, ranging from  $12^\circ$  to  $30^\circ$ , and calculating the frequency of each bin. The tallies are then normalized and provided as part of the observation. This led to much better performance.

Figure 18 shows the returns during training.

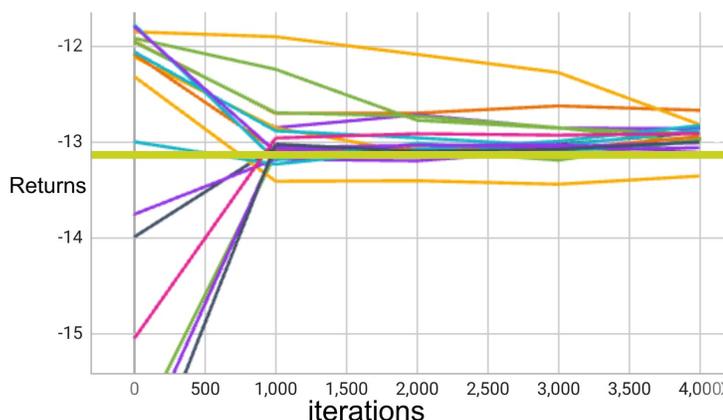


Figure 18: SAC agent Returns of each agent we trained, as well as the baseline in gold, which represents the returns obtained by running the baseline policy currently employed in the real world. As can be clearly seen, most of the agents are able to improve above this policy.

Figure 19 illustrates that the critic, actor, and alpha losses of the various SAC agents converge.

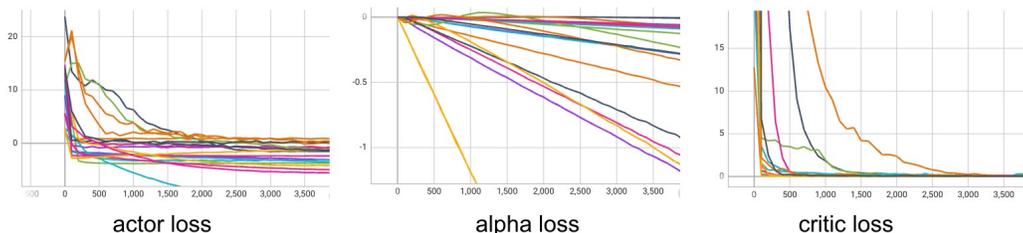


Figure 19: SAC Agent Losses

Our reward function is a weighted, linear combination of the normalized carbon footprint, cost, and comfort levels within the building. While an 8% improvement over the baseline on this scalar reward is significant, we can see the improvements of the SAC agent over the baseline even more clearly when we break down these factors further into physical measures.

For this analysis, we break down the reward into four components that contribute to it, and see how the learned policy compares with the baseline. The components are: setpoint deviation, carbon emissions, electrical energy, and natural gas energy.

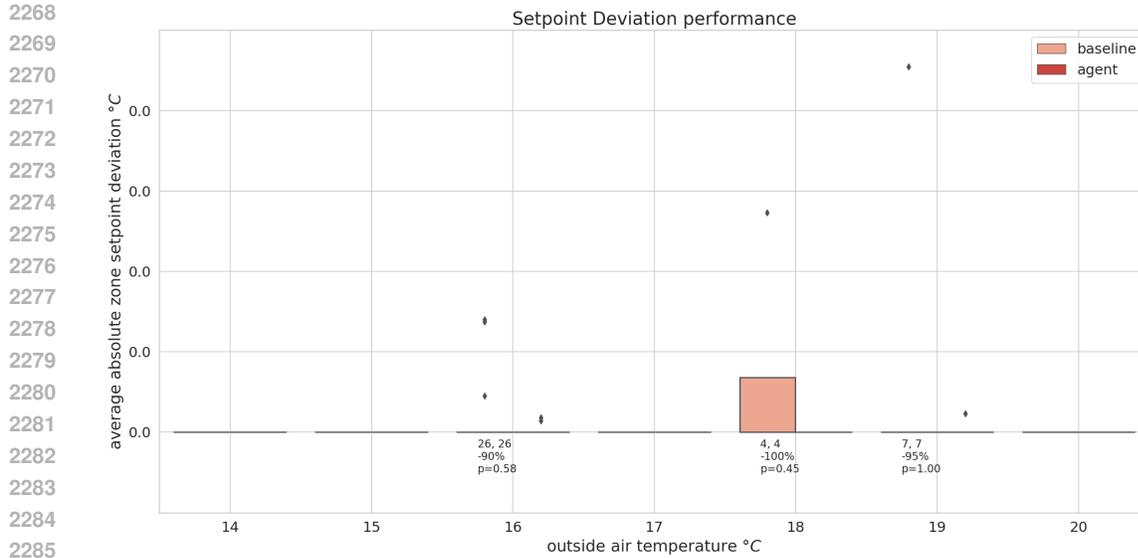


Figure 20: Setpoint Deviation Performance as a function of outside air temperature, which evaluates how well the agent meets comfort conditions compared to the baseline. It is measured as the average number of °C above or below setpoint for all zones in the building. For each outside air degree increment, we include the number of observations for baseline and agent, the percentage change as (baseline - agent) / baseline, and its associated p-score.

Above we display how the baseline and agent compare when it comes to setpoint deviation, the comfort component of the reward function. We show the distribution of deviations grouped by outside air temperatures. While both policies have very minimal setpoint deviation to begin with, the agent strictly improves over the baseline here.

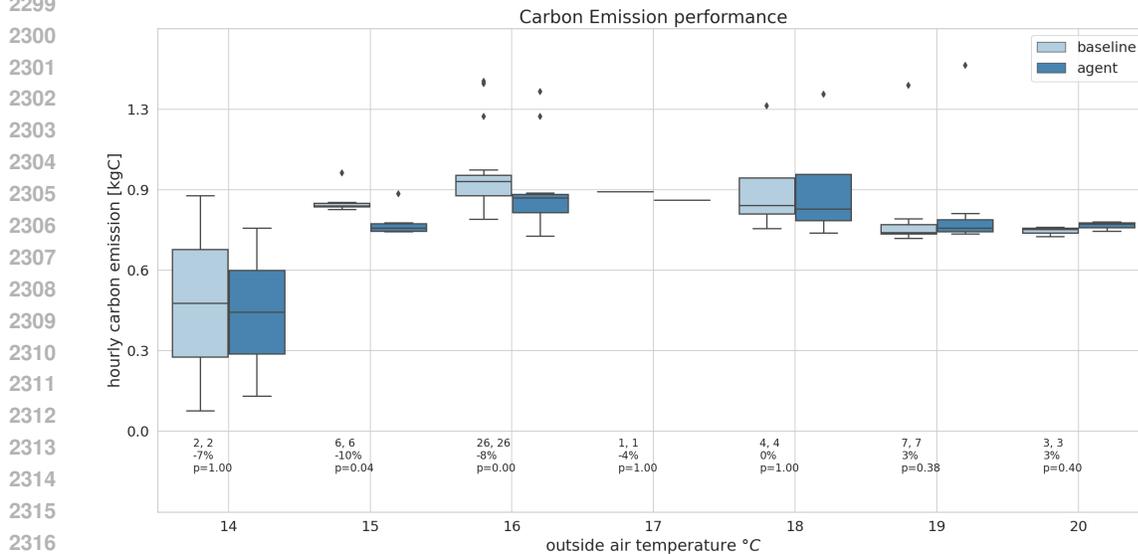
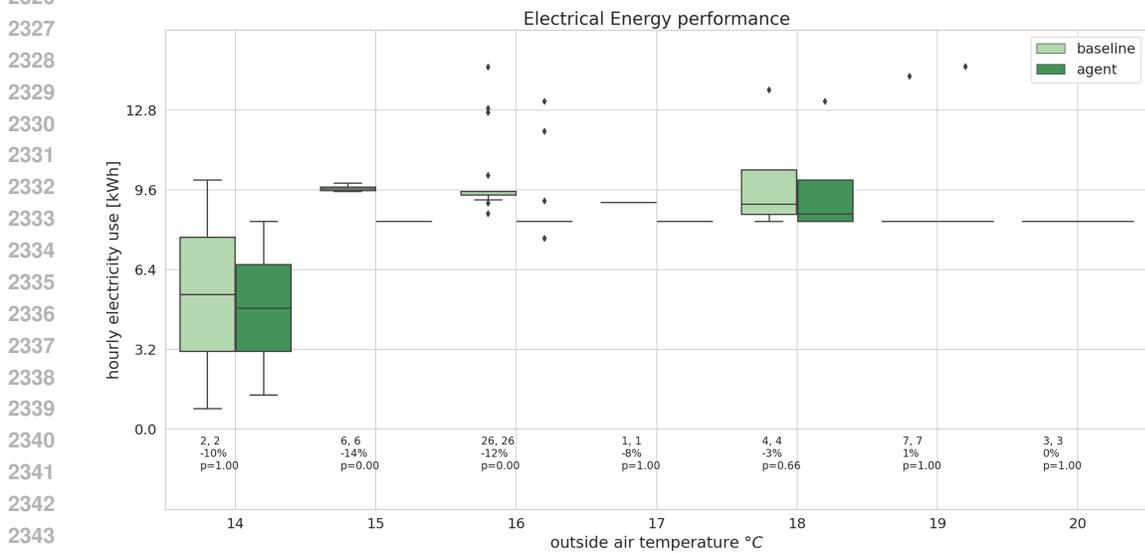


Figure 21: Carbon Emission measures how the agent performs compared to the baseline in terms of the amount of greenhouse gas released from consuming natural gas and electricity. C is combined mass (kgC, or kg Carbon) emitted by non-renewable electricity and natural gas. For each outside air degree increment, we include the number of observations for baseline and agent, the percentage change as (baseline - agent) / baseline, and its associated p-score.

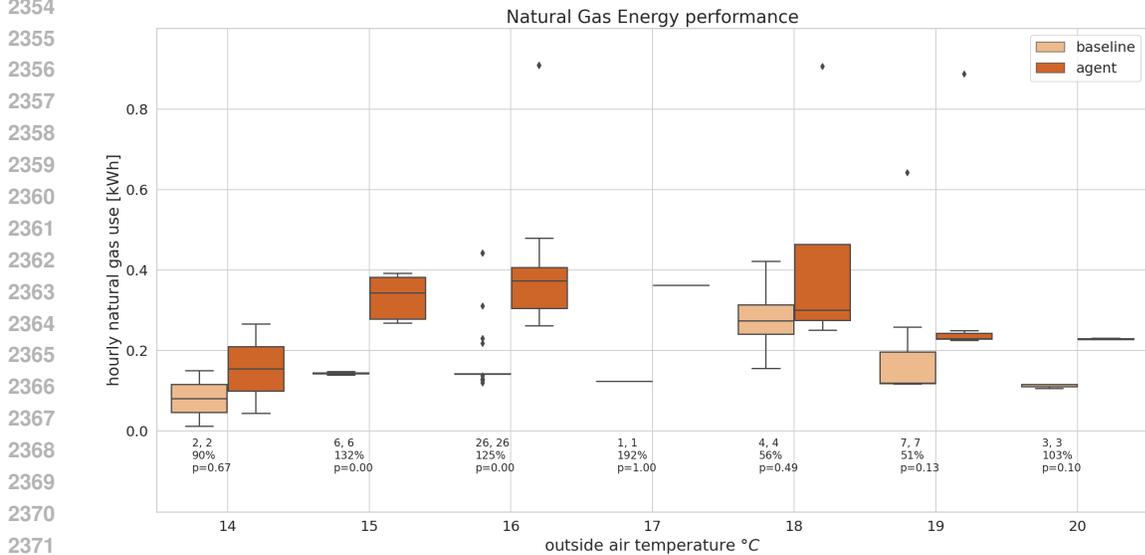
2322 The carbon performance of the agent, as compared with the baseline, is impressive as well. In the  
 2323 temperature range  $14^{\circ}C$  to  $18^{\circ}C$ , the agent is strictly better, and while it is slightly worse for the  
 2324 warmer temperatures, clearly it is a net improvement over the baseline.  
 2325



2345 Figure 22: Electrical Energy Performance measured in energy units (kWh) over a fixed interval  
 2346 for both the agent and the baseline policies. For each outside air degree increment, we include the  
 2347 number of observations for baseline and agent, the percentage change as  $(\text{baseline} - \text{agent}) / \text{baseline}$ ,  
 2348 and its associated p-score.  
 2349

2350

2351 Once again, when it comes to electric performance, the SAC agent is almost strictly better under all  
 2352 temperature ranges.  
 2353



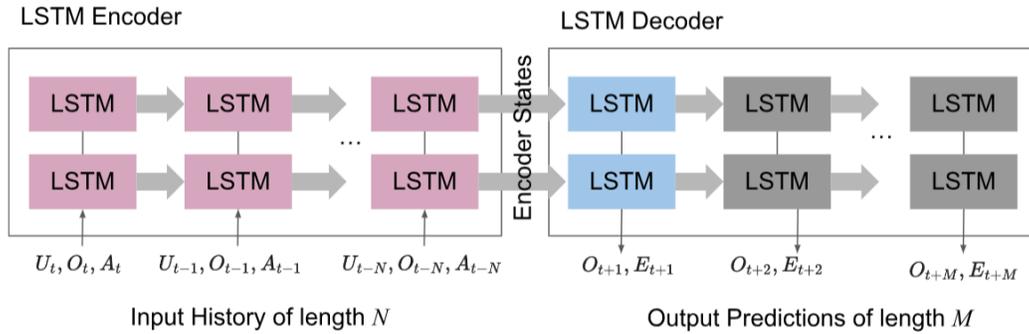
2373 Figure 23: Natural Gas Performance measured in energy units (therm) over a fixed interval for both  
 2374 the agent and the baseline policies. For each outside air degree increment, we include the number  
 2375 of observations for baseline and agent, the percentage change as  $(\text{baseline} - \text{agent}) / \text{baseline}$ , and its  
 associated p-score.

2376 Interestingly, the agent converged on a policy that reduced overall carbon emission while increasing  
 2377 natural gas consumption. This is due to the fact that electricity is generated from non-renewable  
 2378 sources and per unit energy, is significantly more expensive than gas.  
 2379

## 2380 I TRAINING AND EVALUATING A LEARNED DYNAMICS MODEL

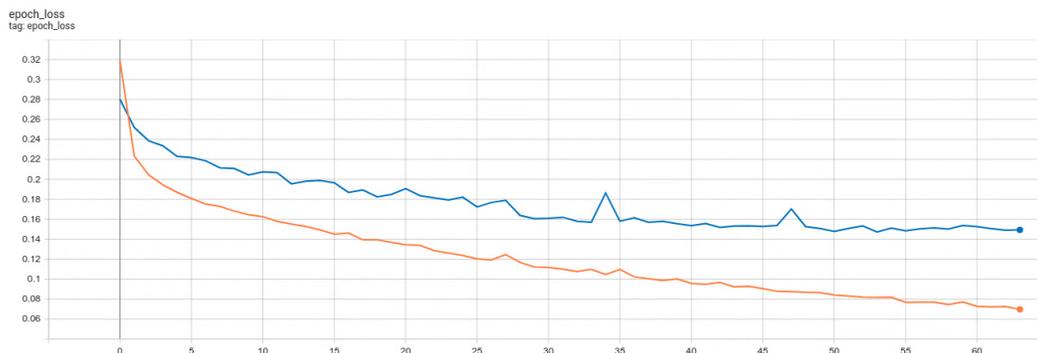
2381  
 2382 Aside from being useful for offline training and for calibrating our simulator, the real world data  
 2383 can also be used to directly learn a regression model that approximates the building dynamics. This  
 2384 model can then be used to train a control agent.  
 2385

2386 As described in the main paper, to demonstrate this approach, building off of earlier work (Velswamy  
 2387 et al., 2017; Sendra-Arranz & Gutiérrez, 2020; Zou et al., 2020; Zhuang et al., 2023), we trained  
 2388 an LSTM to model the building dynamics. We used an encoder-decoder network, where the model  
 2389 takes in a historical sequence of length  $N$  and outputs a prediction sequence of length  $M$ . At each  
 2390 timestep  $t$  in the sequence, the model is given an observation  $O_t$ , action taken by the policy  $A_t$ , and  
 2391 auxiliary state features (such as time of day and weather, that are useful as inputs but need not be  
 2392 predicted)  $U_t$ , and for future timesteps, the model is trained to predict future observations, as well  
 2393 as future reward information (based on predicted energy use and carbon emissions)  $E_t$ . The LSTM  
 2394 model is shown in Figure 24.



2407 Figure 24: Architecture of LSTM building dynamics model

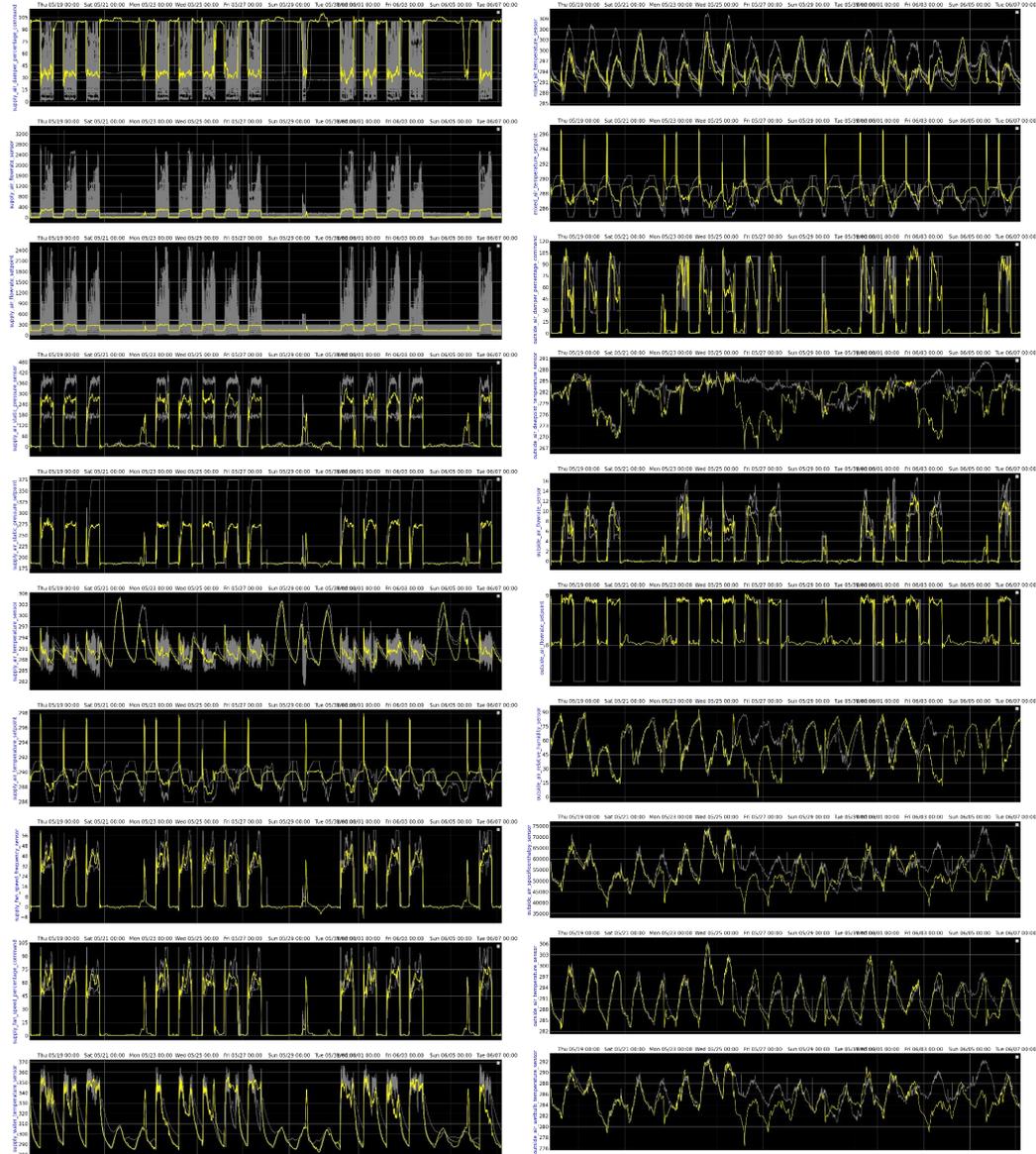
2408  
 2409 We then trained the model to predict the next observation for 65 epochs, plotting training and vali-  
 2410 dation loss, as shown in figure 25



2424 Figure 25: Loss of LSTM building dynamics model, with train loss in orange and validation loss in  
 2425 blue.

2426  
 2427  
 2428 However, loss curves alone do not tell the full story of how well our regression model is reconstruct-  
 2429 ing the signal of the dynamics, so we also included additional evaluations. We had the model predict  
 3 weeks into the future, and then compared the predictions with the ground truth data to ensure the

2430 cyclic patterns of the medians are reproduced. The chart in figure 26 shows 20 measurement time  
 2431 series from the regression models shown in yellow compared to the actual values shown in gray. By  
 2432 inspection, we conclude that the regression building provides good correspondence with the actual  
 2433 real data signals.  
 2434



2435  
 2436  
 2437  
 2438  
 2439  
 2440  
 2441  
 2442  
 2443  
 2444  
 2445  
 2446  
 2447  
 2448  
 2449  
 2450  
 2451  
 2452  
 2453  
 2454  
 2455  
 2456  
 2457  
 2458  
 2459  
 2460  
 2461  
 2462  
 2463  
 2464  
 2465  
 2466  
 2467  
 2468  
 2469  
 2470  
 2471  
 2472  
 2473  
 2474  
 2475  
 2476  
 2477  
 2478  
 2479  
 2480  
 2481  
 2482  
 2483

Figure 26: Detailed analysis of learned dynamics as compared to real data.

## J REAL DATA SAC AGENT TRAINING DETAILS AND PERFORMANCE ANALYSIS

We then trained a SAC agent on the regression environment, much like how we did on the simulator. This gives us a baseline for how to generate a policy purely based on data, without use of the simulator. We used hyper-parameter tuning, and trained 200 agents. The chart in figure 27 shows agent reward progress as the number of trials increased.

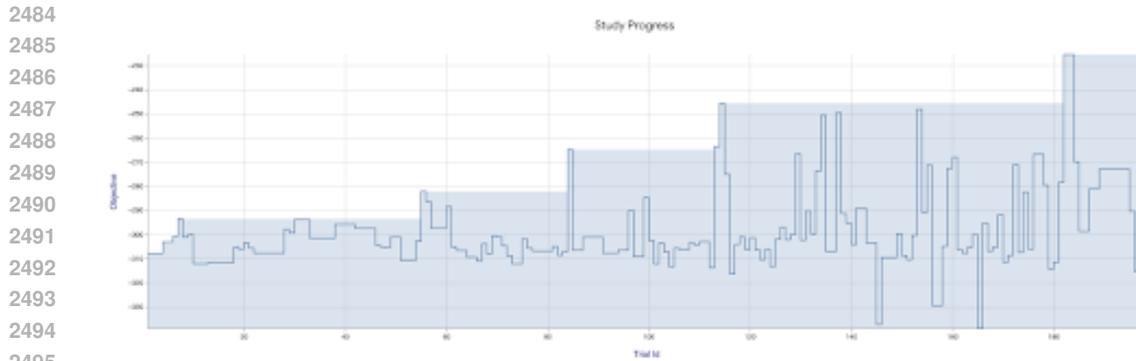


Figure 27: Detailed analysis of learned dynamics as compared to real data.

To compare the learned policy with the baseline, we plotted the two policies in 28. The baseline and agent episode temperature timelines shown below provide a temporal perspective of the environment median zone air temperatures (yellow) and setpoints (white), outside air temperature (blue), and the agent actions on the environment (water temperature setpoints (lime), and air handler temperature setpoints (magenta)). While the regression model under baseline policy correctly represents the weekend setpoint ranges, the regression building applies nearly the weekday setpoint ranges when running under agent control. This is likely due to the agent applying setpoints that regression associates with weekday actions, and incorrectly returns a setpoint that is closer to the weekday. For this reason, we do not evaluate the model’s performance on weekends. Similar to baseline control, the agent ramps up water temperature (lime) at the beginning of the day. However, the agent tends to maintain the water temperature around 80C for substantially longer than baseline control. At first glance, this may seem counterproductive. However, heat exchange is also based on water flow and air flow. Lower supply water temperatures require more airflow to transfer the same amount of heat. Therefore, higher water temperatures do not necessarily result in higher energy consumption. Also, note that the agent does not drop the water temperature as low as the baseline policy, and the agent tends to apply smoother actions compared to the baseline’s rapid oscillation between 40 and 60C. We speculate that one strength of the proposed solution is the agent’s ability to discover better and non-intuitive policies that are unlikely to be chosen by human HVAC technicians. The agent also has a different control policy for the air handlers’ supply air temperatures, shown in magenta. On one air handler’s supply air temperature, the agent tends to operate SB1:AHU:AC 1 at a higher temperature than SB1:AHU:AC 2.

Finally, much like how we did with the simulated agent, we break down the reward into its four components and see how the agent did relative to the baseline on the regression building model.

### J.1 SETPOINT MANAGEMENT PERFORMANCE

The difference in setpoint deviation between agent and baseline was insignificant. However, at 23C the average setpoint deviation was slightly higher, but was still within a narrow window (less than 1/10 C). The setpoint deviation test using the regression model may be slightly optimistic compared to the real building, because the regression model only approximates the zone temperatures with a single median, hiding the larger spread of temperatures throughout the building.

### J.2 CARBON EMISSION PERFORMANCE

In 12 of 19 temperature bins the agent generated less carbon than the baseline. While only two temperature bins (17C, 25C) resulted in confidence greater than 90%, the results indicate a reduction in carbon emission on most of the bins. The agent tends to emit more carbon in the moderate temperature ranges (21, 22C), likely due to a higher setpoint during the day than the baseline. Overall, the agent performs favorably, even though most bins have a low statistical confidence.

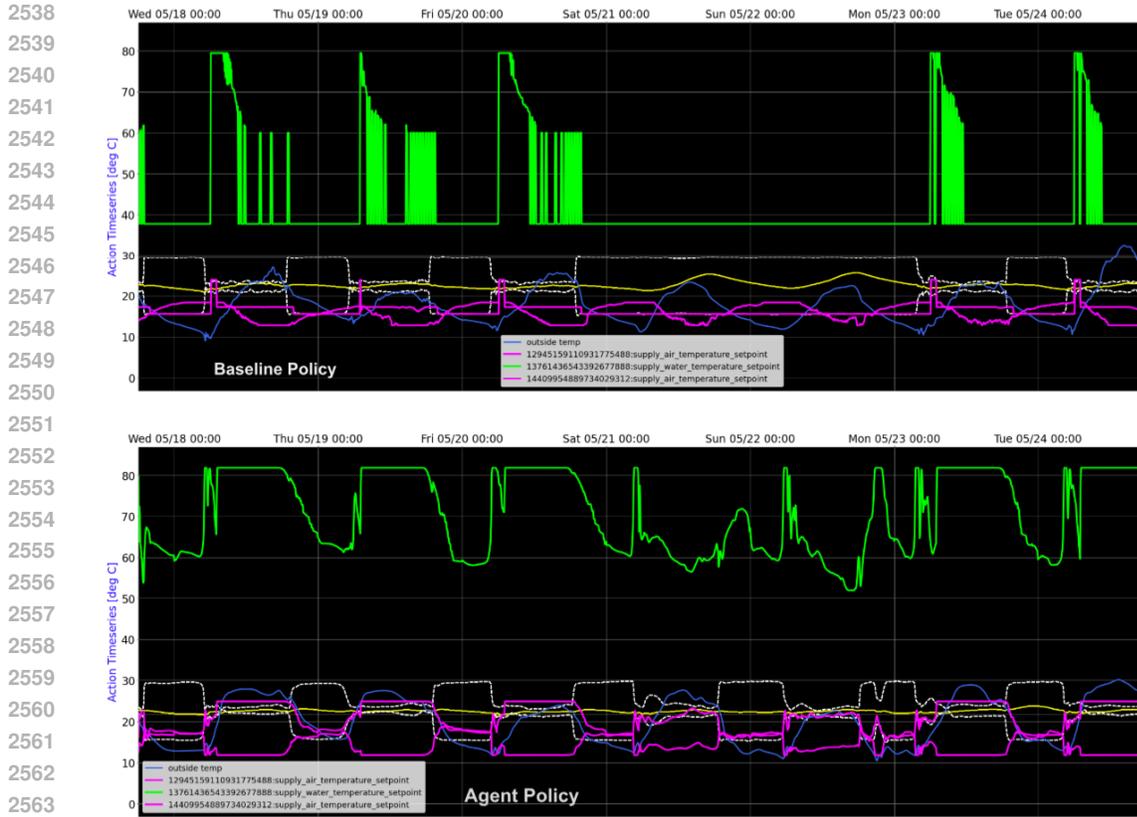


Figure 28: Detailed analysis of learned dynamics as compared to real data.

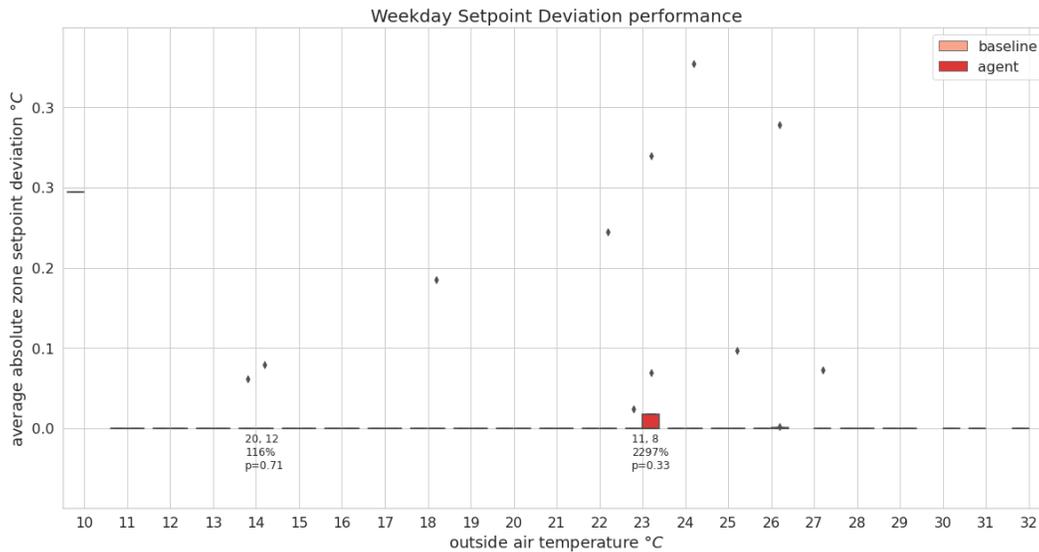


Figure 29: Setpoint Management Performance.

### J.3 ELECTRICAL ENERGY PERFORMANCE

While no temperature bins yielded confidence scores greater than 90%, the agent tends to consume less electricity than the baseline, except for the 21, 22C temperature bins. Under both policies, electricity consumption dramatically increases with outside air temperature.

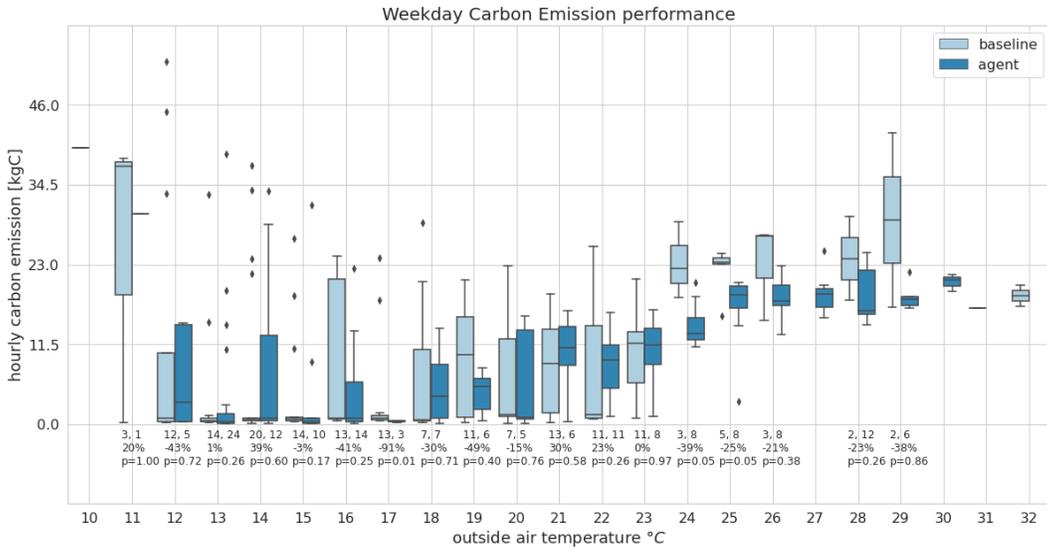


Figure 30: Carbon Emission Performance.

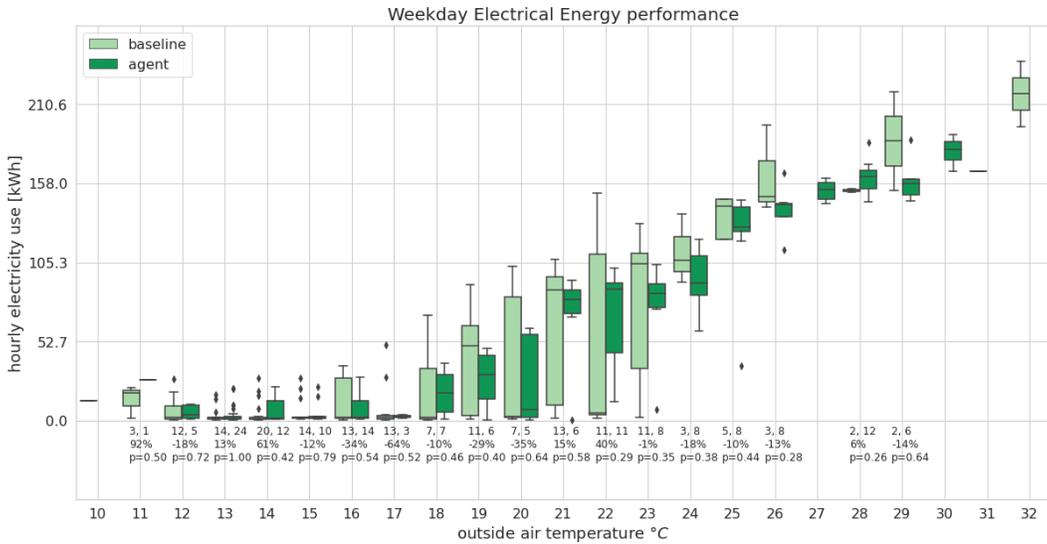


Figure 31: Electrical Energy Performance.

#### J.4 NATURAL GAS ENERGY PERFORMANCE

The agent policy tends to consume less natural gas than the baseline policy, even though only three yielded significant reduction with confidence of at least 90% (17, 24, 25C).

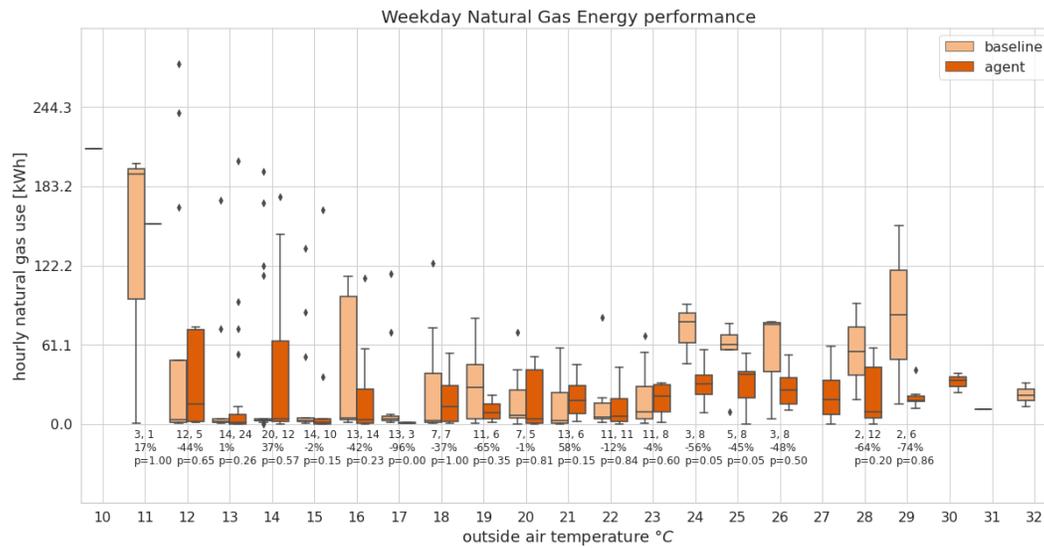


Figure 32: Natural Gas Energy Performance.